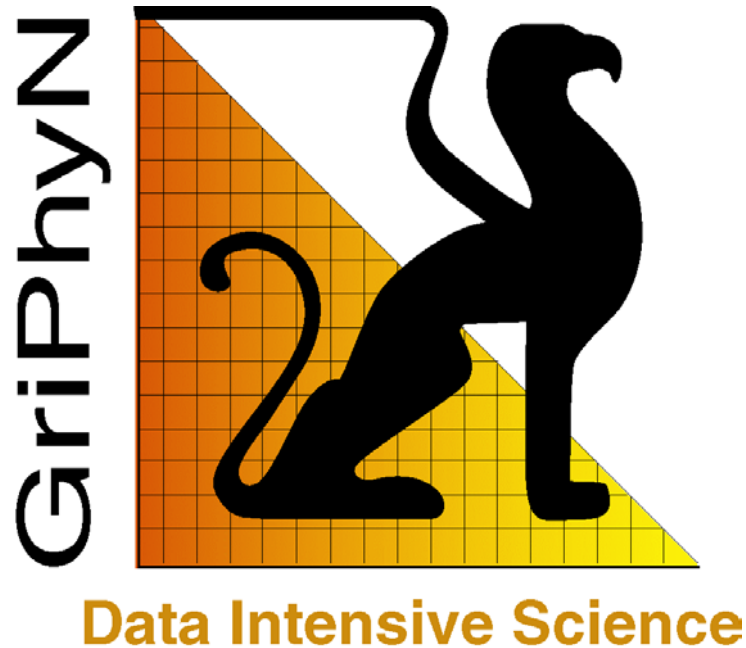


GriPhyN External Advisory Committee Meeting

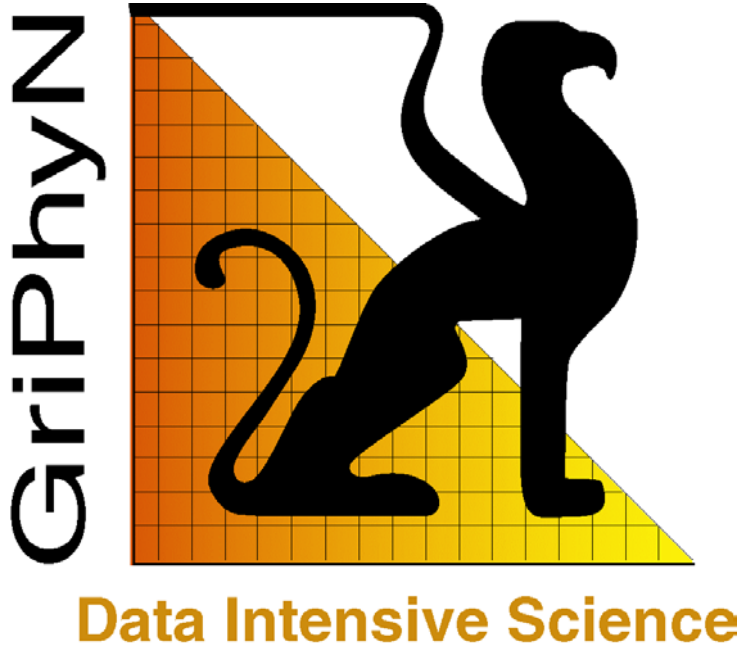


University of Florida
Gainesville, Florida
Jan. 7, 2002

Agenda

Time	Mins	Title	Speaker
8:30	20	Breakfast	
8:50	10	EAC Executive Breakout	
9:00	35+25	GriPhyN Overview	Paul Avery
10:00	30+15	Architecture	Ian Foster
10:45	10+5	Virtual Data Toolkit	Miron Livny
11:00	15	Break	
11:15	30+30	Progress & Planning	Mike Wilde
12:15	20+10	Inter-Project Coordination	Carl Kesselman
12:45	60	Lunch & Discussion	
13:45	15+5	Education and Outreach	Manuela Campanelli
14:05	20+10	CMS Progress	Rick Cavanaugh
14:35	20+10	LIGO Progress	Ewa Deelman
15:05	60	Break + Q&A	
16:05	45	EAC Executive Session	
16:50	20	Report to GriPhyN	
17:10		Adjourn	

GriPhyN Project Overview



Paul Avery
University of Florida
<http://www.phys.ufl.edu/~avery/>
avery@phys.ufl.edu

GriPhyN External Advisory Committee Meeting
Gainesville, Florida
Jan. 7, 2002

Topics

- Overview of Project
- Some observations
- Progress to date
- Management issues
- Budget and hiring
- GriPhyN and iVDGL
- Upcoming meetings
- EAC Issues

Overview of Project

→ GriPhyN basics

- ◆ \$11.9M (NSF) + \$1.6M (matching)
- ◆ 17 universities, SDSC, 3 labs, >40 active participants
- ◆ 4 physics experiments providing frontier challenges

→ GriPhyN funded primarily as an IT research project

- ◆ 2/3 CS + 1/3 physics

→ Must balance and coordinate

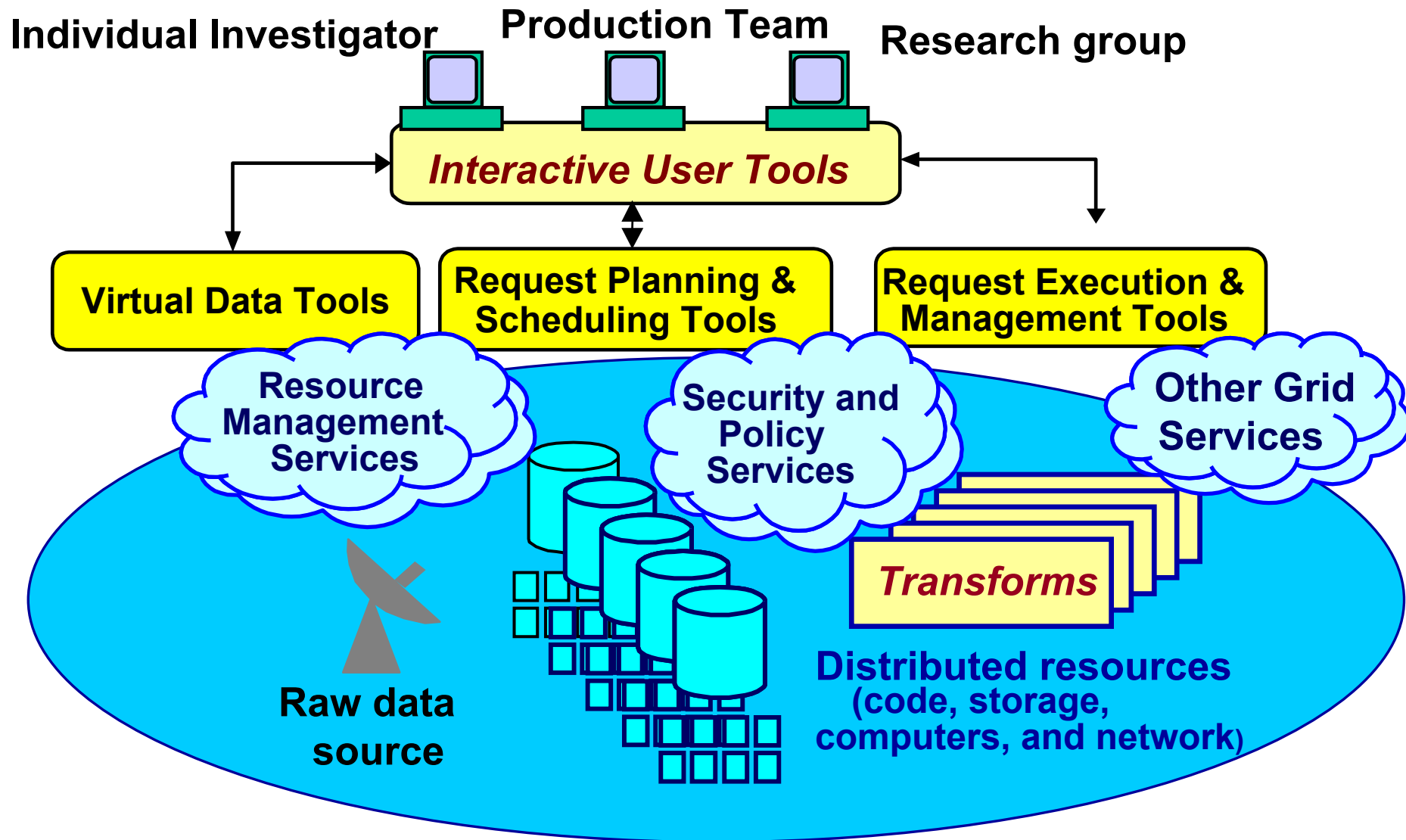
- ◆ Research creativity with project goals and deliverables
- ◆ GriPhyN schedule, priorities and risks with those of 4 experiments
- ◆ Data Grid design and architecture with other Grid projects
- ◆ GriPhyN deliverables with those of other Grid projects



GriPhyN Institutions

- ◆ U Florida
- ◆ U Chicago
- ◆ Boston U
- ◆ Caltech
- ◆ U Wisconsin, Madison
- ◆ USC/ISI
- ◆ Harvard
- ◆ Indiana
- ◆ Johns Hopkins
- ◆ Northwestern
- ◆ Stanford
- ◆ U Illinois at Chicago
- ◆ U Penn
- ◆ U Texas, Brownsville
- ◆ U Wisconsin, Milwaukee
- ◆ UC Berkeley
- ◆ UC San Diego
- ◆ San Diego Supercomputer Center
- ◆ Lawrence Berkeley Lab
- ◆ Argonne
- ◆ Fermilab
- ◆ Brookhaven

GriPhyN: PetaScale Virtual Data Grids



Some Observations

- Progress since April 2001 EAC meeting
 - ◆ Major architecture document draft (with PPDG): Foster talk
 - ◆ Major planning cycle completed: Wilde talk
 - ◆ First Virtual Data Toolkit release in Jan. 2002: Livny talk
 - ◆ Research progress at or ahead of schedule: Wilde talk
 - ◆ Integration with experiments: (SC2001 demos):
Wilde/Cavanaugh/Deelman talks
- Management is "challenging": more later
- Much effort invested in coordination: Kesselman talk
 - ◆ PPDG: shared personnel, many joint projects
 - ◆ EU DataGrid + national projects (US, Europe, UK, ...)
- Hiring almost complete, after slow start
 - ◆ Funds being spent, adjustments being made
- Outreach effort is taking off: Campanelli talk
- iVDGL funded: more later

Progress to Date: Meetings

→ Major GriPhyN meetings

◆ Oct. 2-3, 2000	All-hands	Chicago
◆ Dec. 20, 2000	Architecture	Chicago
◆ Apr. 10-13, 2001	All-hands, EAC	USC/ISI
◆ Aug. 1-2, 2001	Planning	Chicago
◆ Oct. 15-17, 2001	All-hands, iVDGL	USC/ISI
◆ Jan. 7-9, 2002	EAC, Planning, iVDGL	Florida
◆ Feb./Mar. 2002	Outreach	Brownsville
◆ Apr. 24-26 2002	All-hands	Chicago
◆ May/Jun. 2002	Planning	??
◆ Sep. 11-13, 2002	All-hands, EAC	UCSD/SDSC

→ Numerous smaller meetings

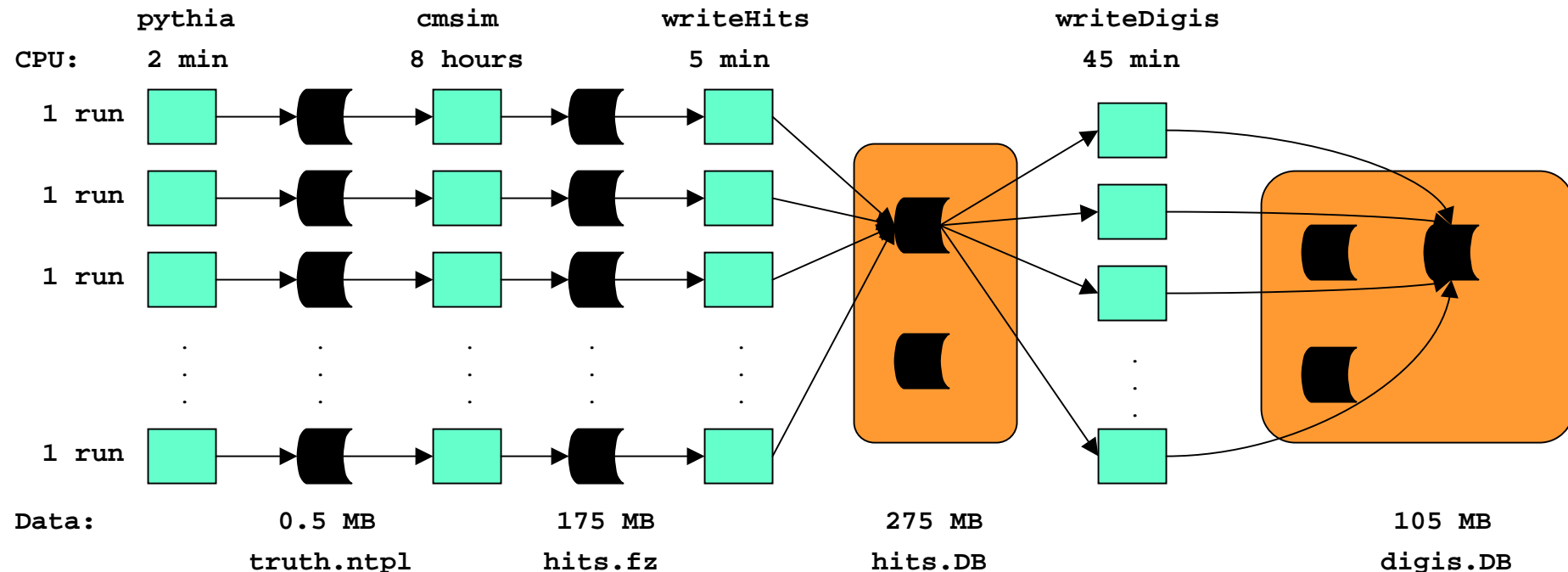
- ◆ CS-experiment
- ◆ CS research
- ◆ Liaisons with PPDG and EU DataGrid

Progress to Date: Conferences

- International HEP computing conference (Sep. 3-7, 2001)
 - ◆ Foster plenary on Grid architecture, GriPhyN
 - ◆ Richard Mount plenary on PPDG, iVDGL
 - ◆ Avery parallel talk on iVDGL
 - ◆ Several talks on GriPhyN research & application work
 - ◆ Several PPDG talks
 - ◆ Grid coordination meeting with EU, CERN, Asia
- SC2001 (Nov. 2001)
 - ◆ Major demos demonstrate integration with expts
 - ◆ LIGO, 2 CMS demos (GriPhyN) + CMS demo (PPDG)
 - ◆ Professionally-made flyer for GriPhyN, PPDG & iVDGL
 - ◆ SC Global (Foster) + International Data Grid Panel (Avery)
- HPDC (July 2001)
 - ◆ Included GriPhyN-related talks + demos
- GGF now features GriPhyN updates

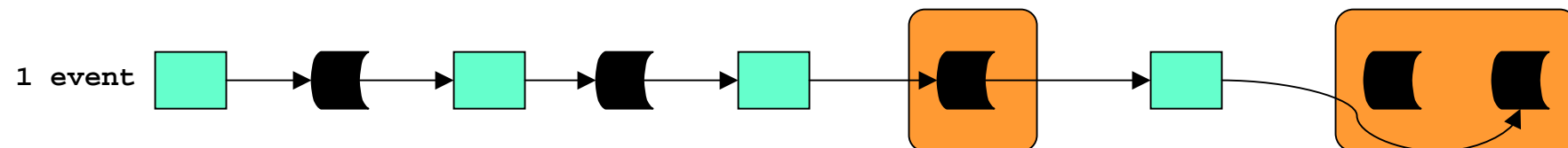


Production Pipeline GriPhyN-CMS Demo

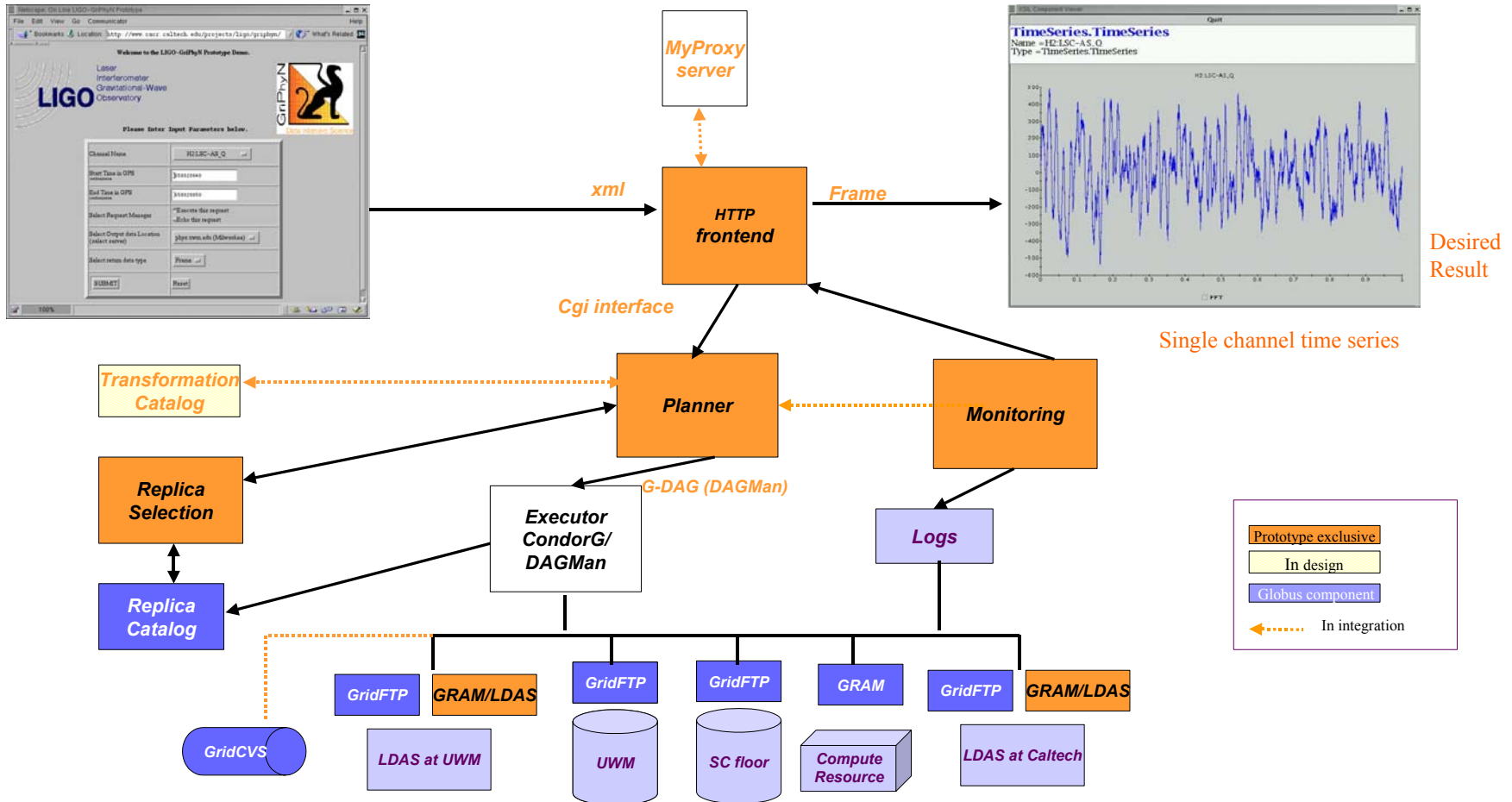


1 run = 500 events

SC2001 Demo Version:

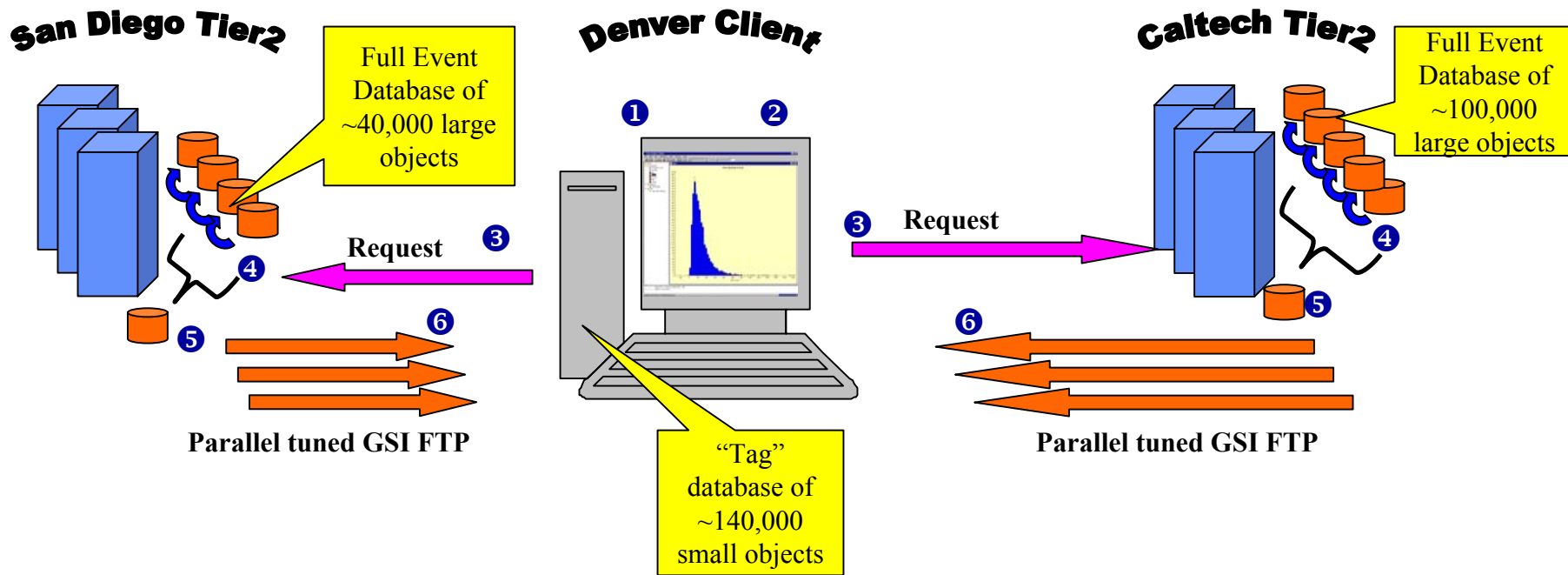


GriPhyN-LIGO SC2001 Demo



GriPhyN CMS SC2001 Demo

http://pcbunn.cacr.caltech.edu/Tier2/Tier2_Overall_JJB.htm



**Bandwidth Greedy Grid-enabled Object Collection Analysis
for Particle Physics**

Management Issues

- Basic challenges from large, dispersed, diverse project
 - ◆ 12 funded institutions + 9 unfunded ones
 - ◆ "Multi-culturalism": CS, 4 experiments
 - ◆ Different priorities and risk equations
- Co-Director concept really works
 - ◆ Share work, responsibilities, blame
 - ◆ CS (Foster) \Leftrightarrow Physics (Avery) \Rightarrow early reality check
 - ◆ Good cop / bad cop useful sometimes
- Project coordinators have helped tremendously
 - ◆ Starting in late July
 - ◆ Mike Wilde Coordinator (Argonne)
 - ◆ Rick Cavanaugh Deputy Coordinator (Florida)

Management Issues (cont.)

→ Overcoming internal coordination challenges

- ◆ Conflicting schedules for meetings
- ◆ Experiments in different stages of software development
- ◆ Joint milestones require negotiation
- ◆ We have overcome these (mostly)

→ Addressing external coordination challenges

- ◆ National: PPDG, iVDGL, TeraGrid, Globus, NSF, SciDAC, ...
- ◆ International: EUDG, LCGP, GGF, HICB, GridPP, ...
- ◆ Networks: Internet2, ESNET, STAR-TAP, STARLIGHT, SURFNet, DataTAG
- ◆ Industry trends: IBM announcement, SUN, ...
- ◆ Highly time dependent
- ◆ Requires lots of travel, meetings, energy

Management Issues (cont.)

→ GriPhyN + PPDG: excellent working relationship

- ◆ GriPhyN: CS research, prototypes
- ◆ PPDG: Deployment
- ◆ Overlapping personnel (particularly Ruth Pordes)
- ◆ Overlapping testbeds
- ◆ Jointly operate iVDGL

→ Areas where we need improvement / advice

- ◆ Reporting system
 - Monthly reports not yet consistent
 - Better carrot? Bigger stick? Better system? Personal contact?
- ◆ Information dissemination
 - Need more information flow from top to bottom
- ◆ Web page
 - Already addressed, changes will be seen soon
- ◆ Using our web site for real collaboration

GriPhyN's Place in the Grid Landscape

→ GriPhyN

- ◆ Virtual data research & infrastructure
- ◆ Advanced planning & execution
- ◆ Fault tolerance
- ◆ Virtual Data Toolkit (VDT)

→ Joint/PPDG

- ◆ Architecture definition
- ◆ Integrating with GDMP and MAGDA
- ◆ Jointly developing replica catalog and reliable data transfer
- ◆ Performance monitoring

→ Joint/others

- ◆ Testbeds, performance monitoring, job languages, ...

→ Needs from others

- ◆ Databases, security, network QoS, ... (Kesselman talk)

→ Focus of meetings over next two months

Management Issues (cont.)

→ Reassessing our work organization (Fig.)

- ◆ ~1 year of experience
- ◆ Rethink breakdown of tasks, responsibilities in light of experience
- ◆ Discussions during next 2 months

→ Exploit iVDGL resources and close connection

- ◆ Testbeds handled by iVDGL
- ◆ Common assistant to iVDGL and GriPhyN Coordinators
- ◆ Common web site development for iVDGL and GriPhyN
- ◆ Common Outreach effort for iVDGL and GriPhyN
- ◆ Additional support from iVDGL for software integration



GriPhyN Management

External Advisory Board

Internet 2
NSF PACIS
DOE Science

Project Directors
Paul Avery
Ian Foster

Project Coordination Group

Project Coordinators
M. Wilde, R. Cavanaugh

System Integration
Carl Kesselman

Industrial Connections
Alex Szalay

Outreach/Education
Manuela Campanelli

Other Grid Projects

Collaboration Board

Physics Experiments

Applications Coord.: H. Newman

ATLAS
(Rob Gardner)

CMS
(Harvey Newman)

LSC(LIGO)
(Bruce Allen)

SDSS
(Alexander Szalay)

VD Toolkit Development Coord.: M. Livny

Requirements Definition & Scheduling
(Miron Livny)

Integration & Testing
(Carl Kesselman – NMI GRIDS Center)

Documentation & Support
(TBD)

CS Research Coord.: I. Foster

Execution Management
(Miron Livny)

Performance Analysis
(Valerie Taylor)

Request Planning & Scheduling
(Carl Kesselman)

Virtual Data
(Reagan Moore)

Technical Coord. Committee Chair: J. Bunn

H. Newman + T. DeFanti
(Networks)

A. Szalay + M. Franklin
(Databases)

R. Moore
(Digital Libraries)

C. Kesselman
(Grids)

P. Galvez + R. Stevens
(Collaborative Systems)



Budget and Hiring

- (Show figures for personnel)
- (Show spending graphs)
- Budget adjustments
 - ◆ Hiring delays give us some extra funds
 - ◆ Fund part time web site person
 - ◆ Fund deputy Project Coordinator
 - ◆ Jointly fund (with iVDGL) assistant to Coordinators

GriPhyN and iVDGL

→ International Virtual Data Grid Laboratory

- ◆ NSF 2001 ITR program \$13.65M + \$2M (matching)
- ◆ Vision: deploy global Grid laboratory (US, EU, Asia, ...)

→ Activities

- ◆ A place to conduct Data Grid tests "at scale"
- ◆ A place to deploy a "concrete" common Grid infrastructure
- ◆ A facility to perform tests and productions for LHC experiments
- ◆ A laboratory for other disciplines to perform Data Grid tests

→ Organization

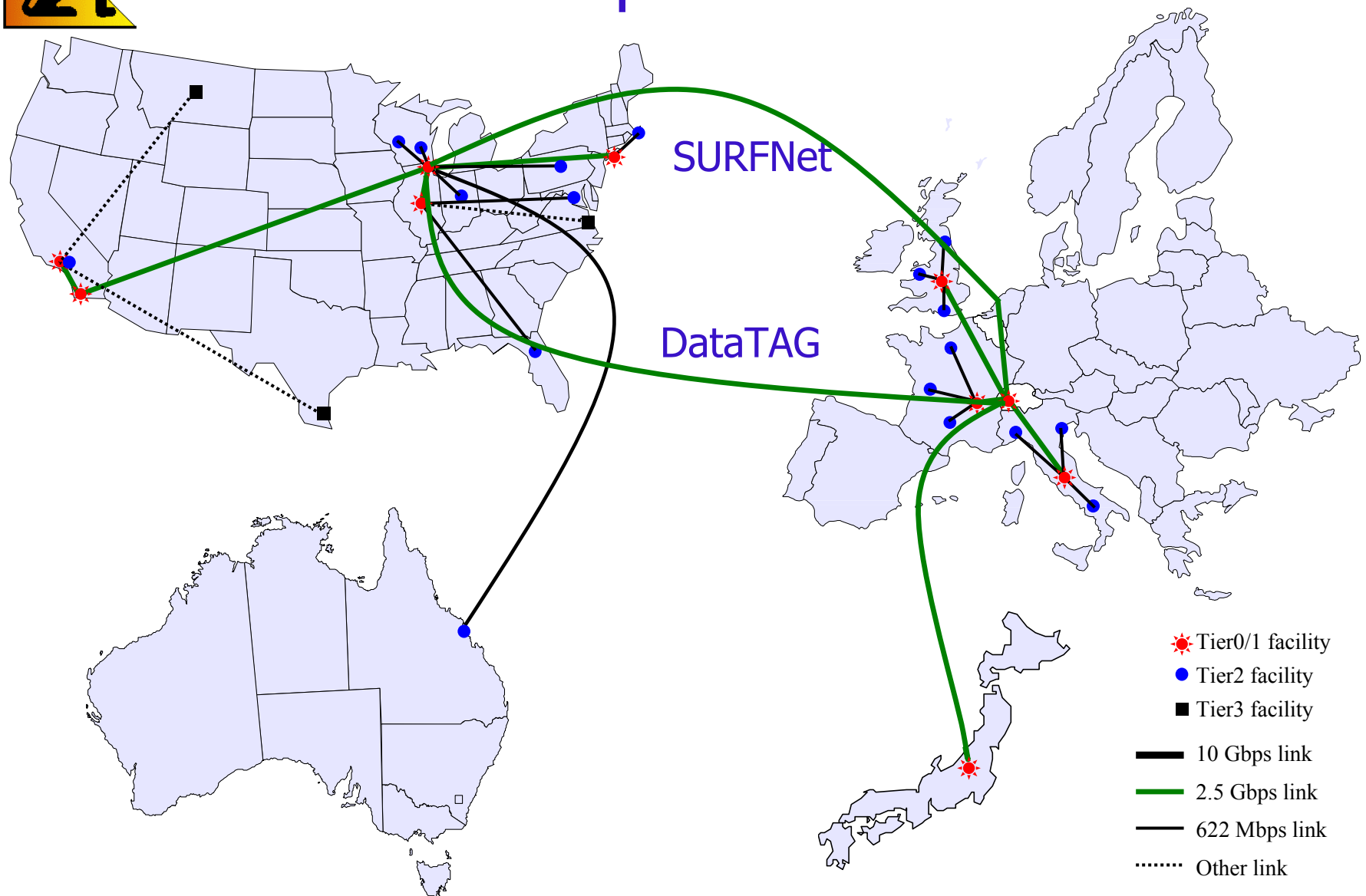
- ◆ GriPhyN + PPDG "joint project"
- ◆ Avery + Foster co-Directors
- ◆ Work teams aimed at deploying hardware/software
- ◆ GriPhyN testbed activities handled by iVDGL



US iVDGL Proposal Participants

◆ U Florida	CMS	}	T2 / Software
◆ Caltech	CMS, LIGO		
◆ UC San Diego	CMS, CS		
◆ Indiana U	ATLAS, iGOC		
◆ Boston U	ATLAS		
◆ U Wisconsin, Milwaukee	LIGO	}	CS support
◆ Penn State	LIGO		
◆ Johns Hopkins	SDSS, NVO		
◆ U Chicago	CS		
◆ U Southern California	CS	}	T3 / Outreach
◆ U Wisconsin, Madison	CS		
◆ Salish Kootenai	Outreach, LIGO		
◆ Hampton U	Outreach, ATLAS		
◆ U Texas, Brownsville	Outreach, LIGO	}	T1 / Labs (not funded)
◆ Fermilab	CMS, SDSS, NVO		
◆ Brookhaven	ATLAS		
◆ Argonne Lab	ATLAS, CS		

iVDGL Map Circa 2002-2003



US-iVDGL Summary Information

→ Principal components (as seen by USA)

- ◆ Tier1 + proto-Tier2 + selected Tier3 sites
- ◆ Fast networks: US, Europe, transatlantic (DataTAG), transpacific?
- ◆ Grid Operations Center (GOC)
- ◆ Computer Science support teams
- ◆ Coordination with other Data Grid projects

→ Experiments

- ◆ HEP: ATLAS, CMS + (ALICE, CMS Heavy Ion, BTeV)
- ◆ Non-HEP: LIGO, SDSS, NVO, biology (small)

→ Proposed international participants

- ◆ 6 Fellows funded by UK for 5 years, work in US
- ◆ US, UK, EU, Japan, Australia (discussions with others)

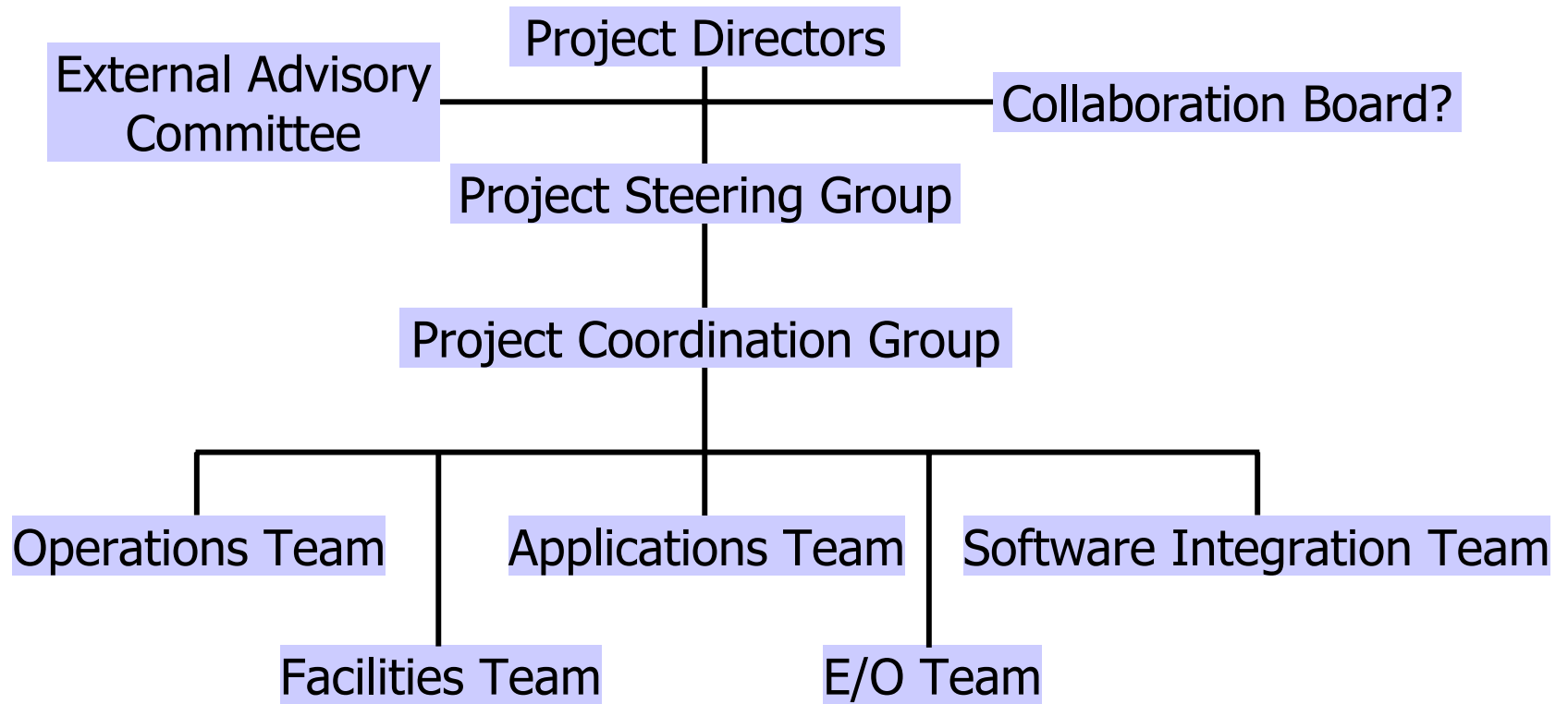


iVDGL Work Team Breakdown

→ Work teams

- ◆ Facilities
- ◆ Software Integration
- ◆ Laboratory Operations
- ◆ Applications
- ◆ Outreach

Initial iVDGL Org Chart





US iVDGL Budget

Units = \$1K

Activity	Year 1	Year 2	Year 3	Year 4	Year 5	Total
CS	338	613	612	756	765	3,084
Outreach	95	166	104	110	115	590
iGOC	115	77	77	232	232	733
ATLAS T2 H/W	157	194	188	58	65	662
ATLAS people	246	338	362	391	392	1,729
CMS T2 H/W	232	192	187	57	65	733
CMS people	234	336	358	390	390	1,708
LIGO T2 H/W	330	58	96	80	48	612
LIGO people	237	358	342	284	286	1,507
SDSS/NVO T2 H/W	160	80	80	40	40	400
SDSS/NVO people	72	88	94	102	102	458
Coordinator	120	180	180	180	180	840
Deputy	50	50	50	50	50	250
EAC costs	20	20	20	20	20	100
Overhead on subcontracts	169	0	0	0	0	169
Reserve	75	0	0	0	0	75
	2,650	2,750	2,750	2,750	2,750	13,650

Upcoming Meetings

→ February / March

- ◆ Outreach meeting, date to be set

→ April 24-26

- ◆ All-hands meeting
- ◆ Perhaps joint with iVDGL

→ May or June

- ◆ Smaller, focused meeting
- ◆ Difficulty setting date so far

→ Sep. 11-13

- ◆ All-hands meeting
- ◆ Sep. 13 as EAC meeting?

EAC Issues

- Can we use common EAC for GriPhyN/iVDGL?
 - ◆ Some members in common
 - ◆ Some specific to GriPhyN
 - ◆ Some specific to iVDGL
- Interactions with EAC: more frequent advice
 - ◆ More frequent updates from GriPhyN / iVDGL?
 - ◆ Phone meetings with Directors and Coordinators?
 - ◆ Other ideas?

GriPhyN/PPDG Virtual Data Grid Architecture

Ian Foster

Mathematics and Computer Science Division
Argonne National Laboratory

Department of Computer Science
The University of Chicago



Overview

- Recap of motivation for architecture effort
- Progress since April 2001
- Architecture as a framework for collaboration
- Web services and databases



Motivation

- We need an architecture so that we can
 - ◆ Coordinate our own activities
 - ◆ Coordinate with other Data Grid projects
 - ◆ Explain to others (experiments, NSF, CS community) what we are doing
- An architecture must:
 - ◆ Facilitate CS research activities by simplifying evaluation of alternatives
 - ◆ Not preclude experimentation with (radically) alternative approaches

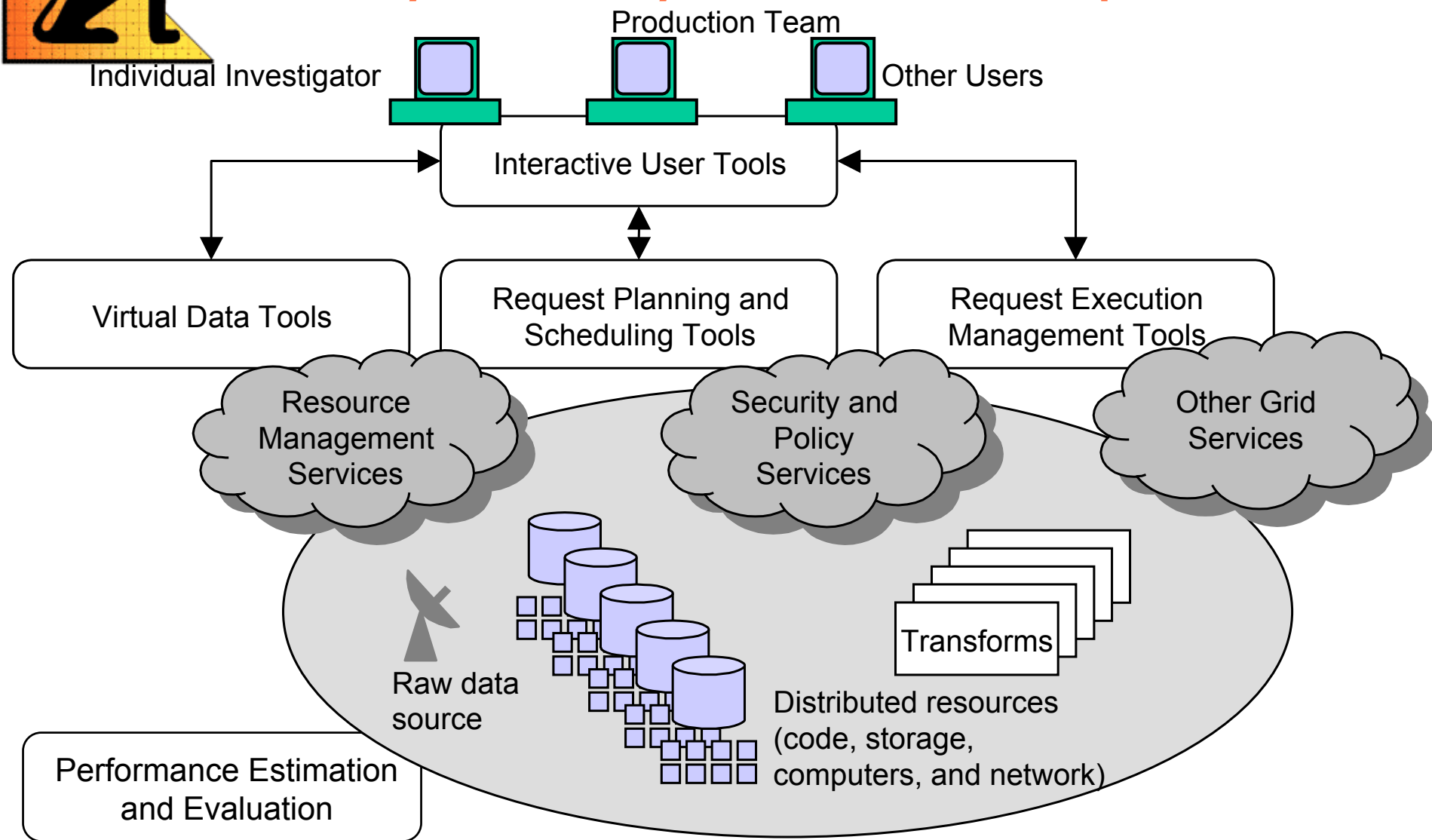


Data Grid Architecture

- Goal: Define requirements and principal elements of a Data Grid system
- Working bottom up, identify elements for which we have solid base of experience and code
- Working top down, identify and prioritize higher-level elements

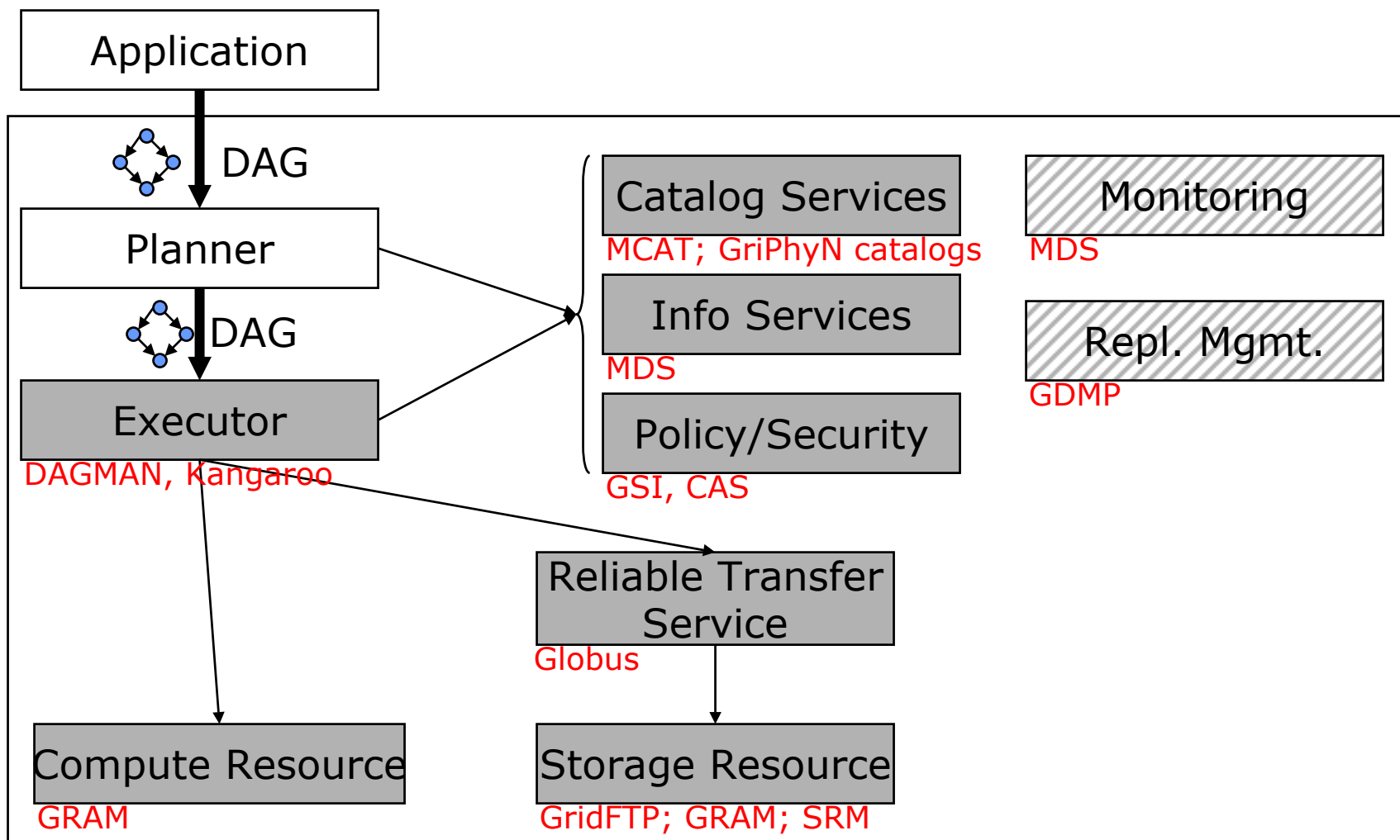


Primary GriPhyN R&D Components





Data Grid Architecture



some [perhaps partial] solution exists

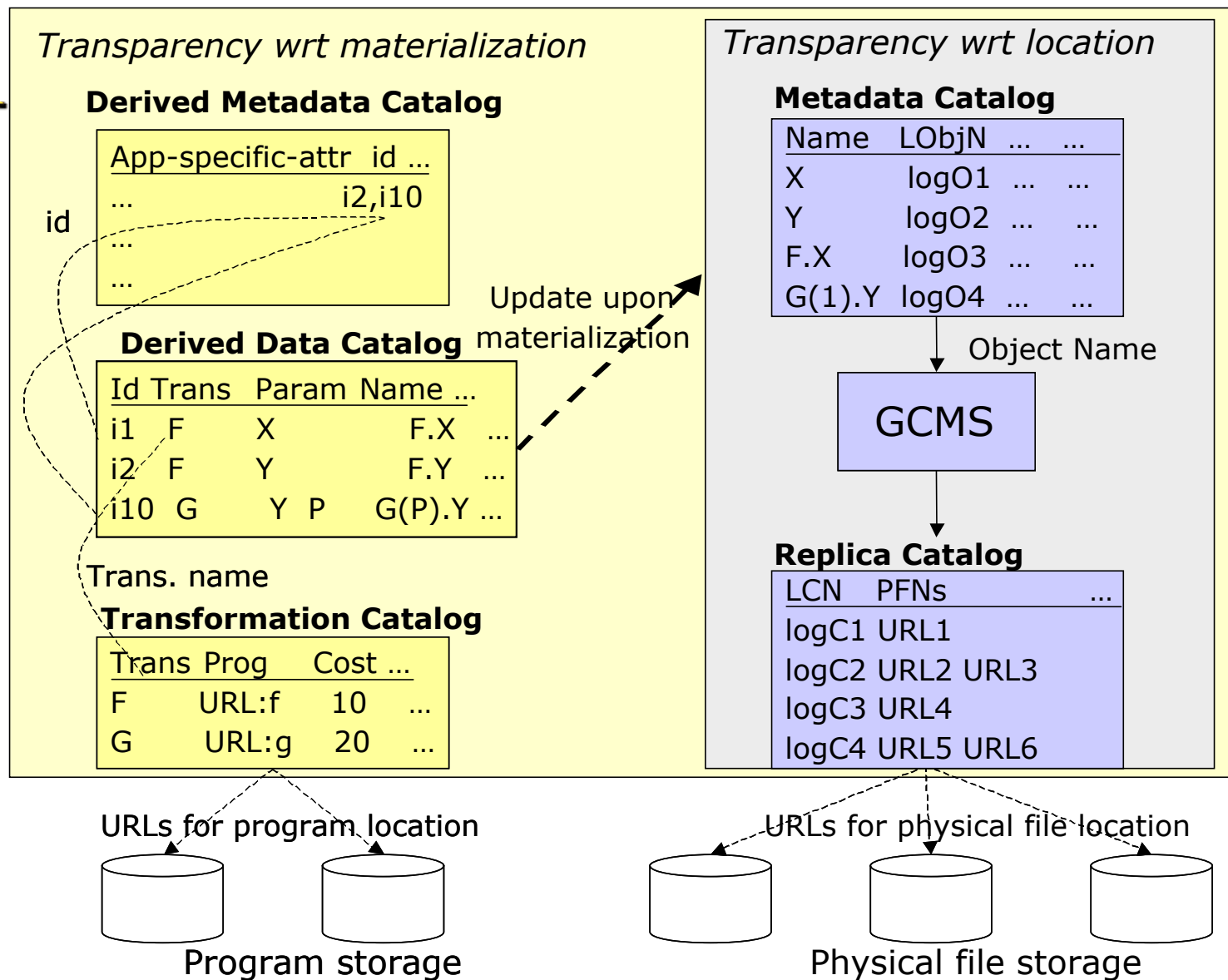
ARGONNE CHICAGO



Progress Since April 2001

- Major architecture revisions and extensions
 - ◆ E.g., virtual data catalog & scheduling tools
- Establishment of collaboration with PPDG
 - ◆ Architecture document now joint
- Adoption by expts of DGA/VDT components
 - ◆ E.g., DAGMAN, GridFTP
- New projects instantiating missing pieces
 - ◆ E.g., Giggle RLS, DAGMAN extensions, VDL
- v2.0 released; significant progress on v3.0

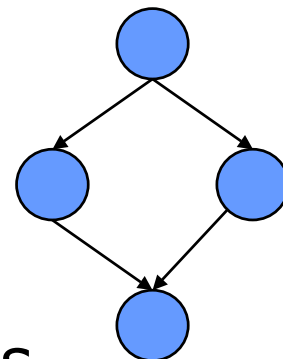
Catalog Architecture





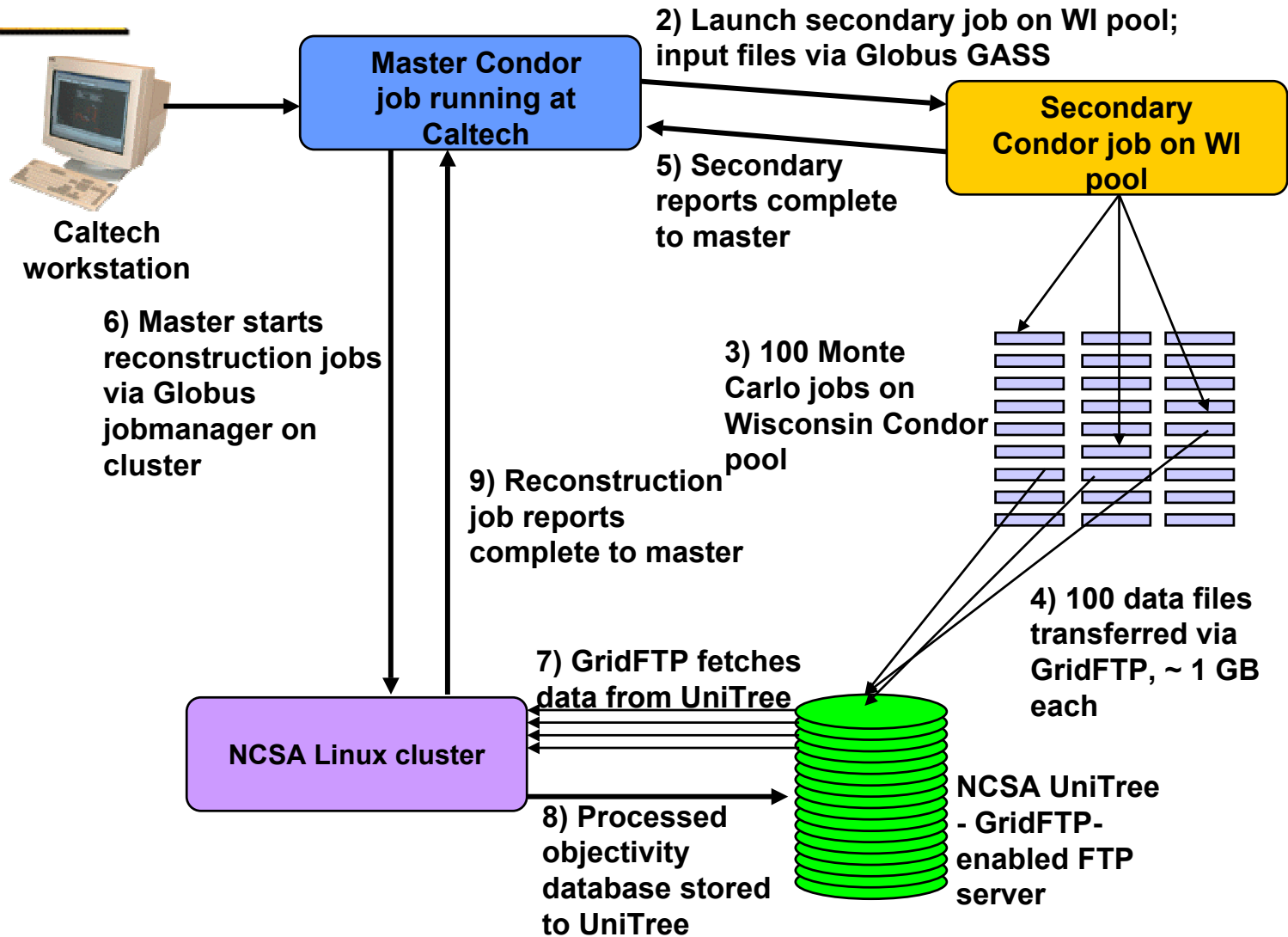
Executor

- Requirements
 - ◆ Reliable management of the execution of a set of computations and data movements
- Current status
 - ◆ UW DAGMan, which executes a supplied directed acyclic graph (DAG)
- Future directions
 - ◆ Error handling and recovery
 - ◆ Ability to express alternative strategies
 - ◆ Beyond DAGs

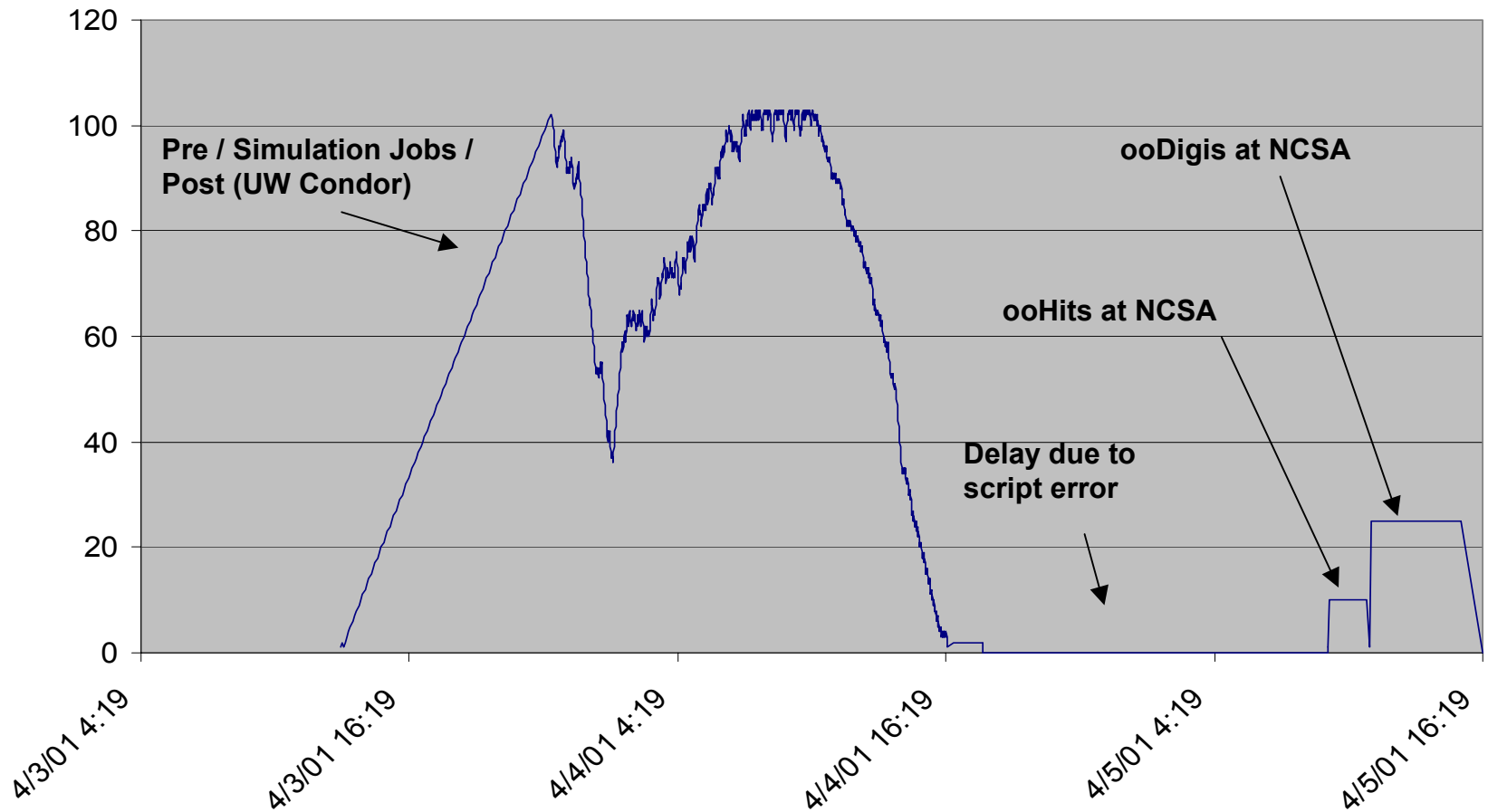




Adoption by Experiments: E.g., CMS Data Reconstruction



Trace of a Condor-G Physics Run





New Projects; Architecture as a Roadmap for Collaboration

- “Giggle” replica location service
 - ◆ Joint design, implementation, & evaluation effort with CERN EU DataGrid team
- Replica management, other higher-level data management tools
 - ◆ Joint work with elements of PPDG and EU DataGrid
- Monitoring and discovery requirements
 - ◆ Joint with PPDG, iVDGL
- DAGMAN extensions and next steps



Web Services and Databases

- Q: How does DGA relate to
 - ◆ Web services?
 - WSDL as interface definition language, SOAP as transport
 - Other services for discovery, workflow, etc.
 - ◆ Databases?
 - Importance for HEP a topic of debate, but clearly critical for many disciplines
 - ◆ In both cases, significant industrial adoption
- A: We have a plan, and partnerships in place to execute this plan during 2002
 - ◆ Only partially GriPhyN/DGA-specific



The Basis for the Plan: Open Grid Services Architecture

- Grid Service = Web service with specified codified behaviors and interfaces, e.g.
 - ◆ Global naming; reliable, secure creation/mgmt of transient, stateful remote service instances
 - ◆ Standard factory, discovery, etc., interfaces
- Roll-out planned during 2002
 - ◆ Evolution of existing protocols, e.g.
 - GRAM-2: Remote management protocol
 - GSI: Restricted delegation, etc.
 - ◆ New protocols & services: e.g., databases



OGSA and DGA

- We are relying on various OGSA components for future VDT components
 - ◆ E.g., GRAM-2 for reservation, restricted delegation for security
 - ◆ These will arrive during 2002
- In the meantime, we are starting to package selected services as Grid services
 - ◆ E.g., reliable data transfer (with LBNL), replica location (with CERN), etc.

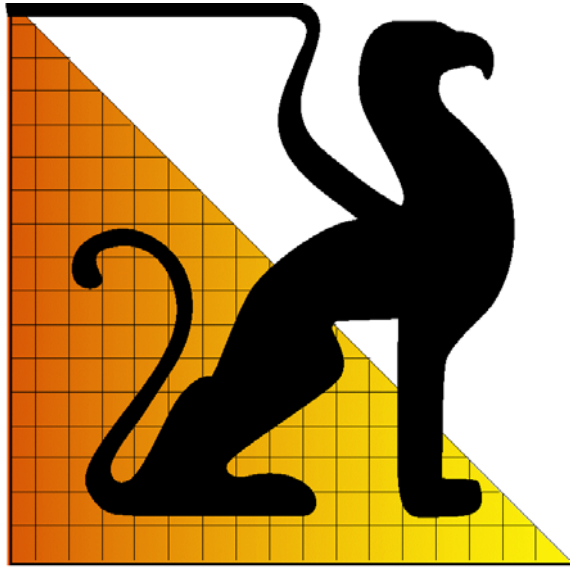


Next Steps

- V3 to be released first quarter 2002
 - ◆ After extensive review by GriPhyN, PPDG
 - ◆ Integrate all latest thinking
- Use this
 - ◆ To drive R&D for 2002-2003
 - ◆ As a basis for discussions with EU DataGrid, LHC Grid, and TeraGrid projects
- Initiate new efforts
 - ◆ Discovery and monitoring
 - ◆ Virtual data catalogs and tools

VDT Status Report

GriPhyN



Data Intensive Science

Miron Livny
Computer Sciences Department
University of Wisconsin-Madison
miron@cs.wisc.edu

GriPhyN External Advisory Committee Meeting
Gainesville, Florida
Jan. 7, 2002

The Mission

“As discussed in Section C.1.d, a primary **GriPhyN** deliverable will be a suite of *virtual data services* and *virtual data tools* designed to support a wide range of applications. The development of this *Virtual Data Toolkit (VDT)* will enable the real-life experimentation needed to evaluate **GriPhyN** technologies. The VDT will also serve as a primary technology transfer mechanism to the four physics experiments and to the broader scientific community”.

“VDT development will be led by Livny, under the direction of the **GriPhyN** Technical Committee (see below).”

VDT 1.0

→ The first version of VDT was defined to include the following components:

◆ VDT- Server

- Condor (version 6.3.1) - Local cluster management and scheduling
- GDMP (version 2.0 beta) - File replication/mirroring.
- Globus Toolkit (version 2.0 beta) - GSI, GRAM, MDS, GridFTP, Replica Catalog & Management all packaged with GPT.

◆ VDT - Client

- Condor-G (version 6.3.1) - Local management of Grid jobs.
- DAGMan - Support Directed Acyclic Graphs (DAGs) of Grid jobs.
- Globus Toolkit (version 2.0 beta) - Client side of GSI, GRAM, GridFTP & Replica Catalog & Management all packaged with GPT.

◆ VDT - Developer

- ClassAd (version 1.0) - Supports collections and Matchmaking
- Globus Toolkit (version 2.0) - Grid APIs

Status

- VDT Developers recruited, hired and organized (1.5 at UC and 2 at UW)
- Components of VDT 1.0 identified, modified, packaged and tested
- Early versions of VDT 1.0 used to “power” SC’01 demos delivering “end-to-end” functionality to real-life applications
- VDT web page constructed.

Next Steps

- Develop and implement a packaging and installation framework for VDT (to be completed by 2Q FY02).
- Establish support procedures and infrastructure.
- Formalize relationship with "external" sources of VDT components.
- Develop requirements and procedures for the timing and content of new releases (focus of our meeting tomorrow).
- Develop scenarios for SC'02 demos/challenges and link them to VDT development and release schedules

VDT 2.0 (3Q FY02)

- Virtual Data Catalog structures and VDL engine
 - VDL and rudimentary centralized planner / executor
- Community Authorization Server
 - Initial Grid Policy Language
- The Network Storage (NeST) appliance
- User login management tools
- A Data Placement (DaP) job manager

Questions and Comments ...

GriPhyN Status and Project Plan

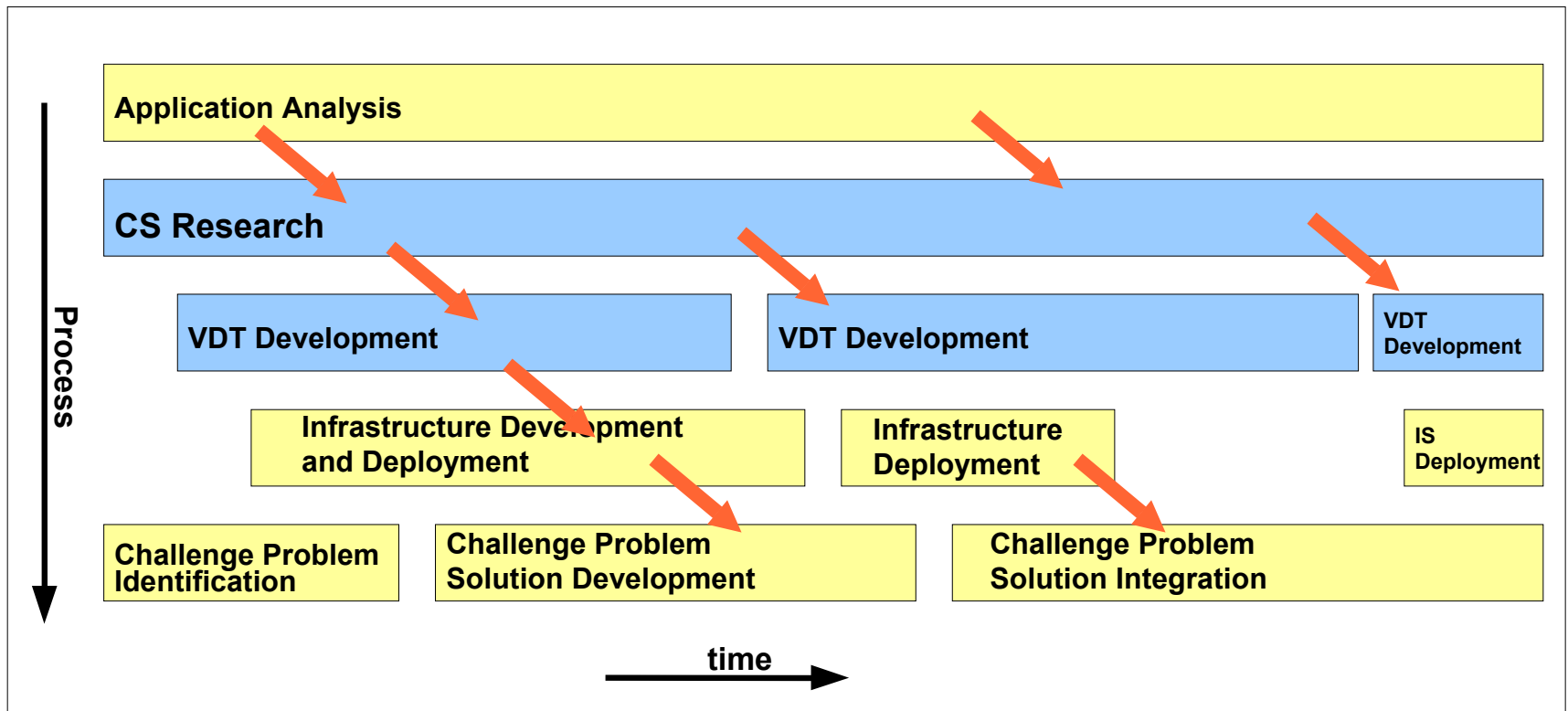
Mike Wilde

Mathematics and Computer Science Division
Argonne National Laboratory

Planning Goals

- Clarify our vision and direction
 - ◆ Know how we can make a difference in science and computing!
- Map that vision to each application
 - ◆ Create concrete realizations of our vision
- Organize as cooperative subteams with specific missions and defined points of interaction
- Coordinate our research programs
- Shape toolkit to meet challenge-problem needs
- “Stop, Look, and Listen” to each experiment’s need
 - ◆ Excite the customer with our vision
 - ◆ Balance the promotion of our ideas with a solid understanding of the size and nature of the problems

Project Approach



Project Activities

- Research

- ◆ Experiment Analysis

- Use cases, statistics, distributions, data flow patterns, tools, data types, HIPO

- ◆ Vision Refinement

- ◆ Attacking the “hard problems”

- Virtual data identification and manipulation
 - Advanced resource allocation and execution planning
 - Scaling this up to Petascale

- ◆ Architectural Refinement

- Toolkit Development

- Integration

- ◆ Identify and Address Challenge Problems
 - ◆ Testbed construction

- Support

- Evaluation

Research Milestone Highlights

- Y1: Execution framework
Virtual data prototypes
- Y2: Virtual data catalog w/glue language
Integ w/ scalable replica catalog service
Initial resource usage policy language
- Y3: Advanced planning, fault recovery
Intelligent catalog
Advanced policy languages
- Y4: Knowledge management and location
- Y5: Transparency and usability
Scalability and manageability

Research Leadership Centers

- Virtual Data:
 - ◆ Chicago (VDC, VDL, KR), ISI (Schema)
 - ◆ Wisconsin (NeST), SDSC (MCAT,SRB)
- Request Planning
 - ◆ ISI (algorithms), Chicago (policy), Berkeley (query optimization)
- Request Execution
 - ◆ Wisconsin
- Fault Tolerance
 - ◆ SDSC
- Monitoring
 - ◆ Northwestern
- User Interface
 - ◆ Indiana

Project Status Overview

- Year 1 research fruitful
 - ◆ Virtual data, planning, execution, integration—demonstrated at SC2001
- Research efforts launched
 - ◆ 80% focused – 20% exploratory
- VDT effort staffed and launched
 - ◆ Yearly major release; VDT1 close; VDT2 planned; VDT3-5 envisioned
- Year 2 experiment integrations high level plans done; detailed planning underway
- Long term vision refined and unified

Milestones: Architecture

- Early 2002:
 - ◆ Specify interfaces for new GriPhyN functional modules
 - Request Planner
 - Virtual Data Catalog service
 - Monitoring service
 - ◆ Define how we will connect and integrate our solutions, e.g.:
 - Virtual data language
 - Multiple-catalog integration
 - DAGman graphs
 - Policy language
 - CAS interaction for policy lookup and enforcement
- Year-end 2002: phased migration to a web-services based architecture

Status: Virtual Data

- Virtual Data
 - ◆ First version of a catalog structure built
 - ◆ Integration language “VDL” developed
 - ◆ Detailed transformation model designed
- Replica location service at Chicago & ISI
 - ◆ Highly scalable and fault tolerant
 - ◆ Soft-state distributed architecture
- NeSt at UW
 - ◆ Storage appliance for Condor
 - ◆ Treats data transfer as a job step

Milestones: Virtual Data

- Year 2:
 - ◆ Local Virtual Data Catalog Structures (relational)
 - ◆ Catalog manipulation language (VDL)
 - ◆ Linkage to application metadata
- Year 3: Handling multi-modal virtual data
 - ◆ Distributed virtual data catalogs (based on RLS)
 - ◆ Advanced transformation signatures
 - ◆ Flat, objects, OODBs, relational
 - ◆ Cross-modal dependency tracking
- Year 4: Knowledge representation
 - ◆ Ontologies; data generation paradigms
 - ◆ Fuzzy dependencies and data equivalence
- Year 5: Finalize Scalability and Manageability

Status: Planning and Execution

- Planning and Execution
 - ◆ Major strides in execution environment made with Condor, CondorG, and DAGman
 - ◆ DAGs evolving as pervasive job specification model with the virtual data grid
 - ◆ Large-scale CMS production demonstrated on 3-site wide-area multi-organization grid
 - ◆ LIGO demonstrated full GriPhyN integration
 - ◆ Sophisticated policy language for grid-wide resource sharing under design at Chicago
 - ◆ Knowledge representation research underway at Chicago
 - ◆ Research in ClassAds explored in Globus context
- Master/worker fault tolerance at UCSD
 - ◆ Design proposed to extend fault tolerance of Condor masters

Milestones: Request Planning

- Year 2:
 - ◆ Prototype planner as a grid service module
 - ◆ Initial CAS and Policy Language Integration
 - ◆ Refinement of DAG language with data flow info
- Year 3:
 - ◆ Policy enhancements: dynamic replanning (based on Grid monitoring), cost alternatives and optimizations
- Year 4:
 - ◆ Global planning with policy constraints
- Year 5:
 - ◆ Incremental global planning
 - ◆ Algorithms evaluated, tuned w/ large-scale simulations

Milestones: Request Execution

- Year 2:
 - ◆ Request Planning and Execution
 - Striving for increasingly greater resource leverage with increasing both power AND transparency
 - Fault tolerance – keeping it all running!
 - ◆ Initial CAS and Policy Language Integration
 - ◆ Refinement of DAG language with data flow info
 - ◆ Resource utilization monitoring to drive planner
- Year 3:
 - ◆ Resource co-allocation with recovery
 - ◆ Fault tolerant execution engines
- Year 4:
 - ◆ Execution adapts to grid resource availability changes
- Year 5:
 - ◆ Simulation-based algorithm eval and tuning

Status: Supporting Research

- Joint PPDG-GriPhyN Monitoring group
 - ◆ Meeting regularly
 - ◆ Use-case development underway
- Research into monitoring, measurement, profiling, and performance predication
 - ◆ Underway at NU and ANL
- GRIPE facility for Grid-wide user and host certificate and login management
- GRAPPA portal for end-user science access

Status – Experiments

- ATLAS

- ◆ 8-site testgrid in place
- ◆ data and metadata management prototypes evolving
- ◆ Ambitious Year-2 plan well refined – will use numerous GriPhyN deliverables

- CMS

- ◆ Working prototypes of production and distributed analysis, both with virtual data
- ◆ Year-2 plan – simulation production – underway

- LIGO

- ◆ Working prototypes of full VDG demonstrated
- ◆ Year-2 plan well refined and development underway

- SDSS

- ◆ Year-2 plan well refined
- ◆ Challenge problem development underway
- ◆ close collaboration with Chicago on VDC

Year 2 Plan: ATLAS

- ATLAS-GriPhyN Challenge Problem I
 - ◆ ATLAS DC0: 10M events, $O(1000)$ CPUs
 - ◆ Integration of VDT to provide uniform distributed data access
 - ◆ Use of GRAPPA portal, possibly over DAGman
 - ◆ Demo ATLAS SW Week – March 2002

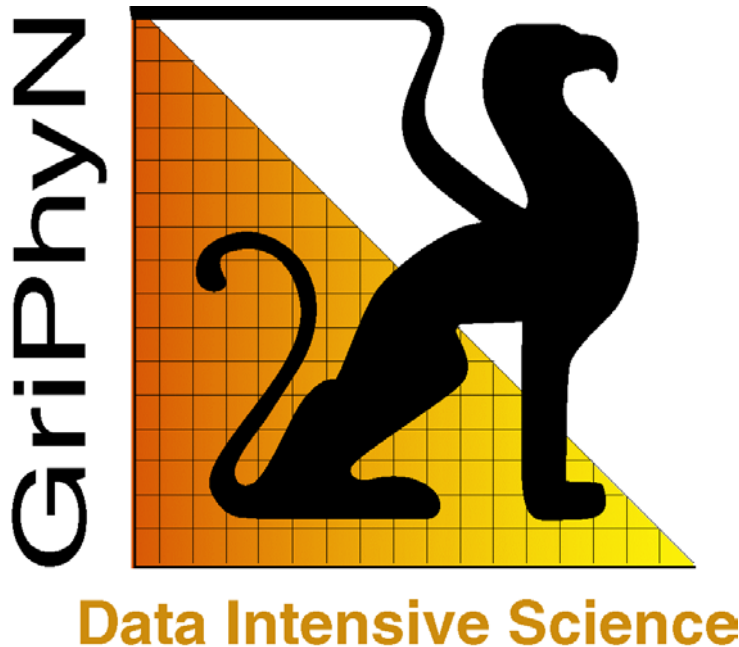
Year 2 Plan: ATLAS

- ATLAS-GriPhyN Challenge Problem II
 - ◆ Virtualization of pipelines to deliver analysis data products: reconstructions and metadata tags
 - ◆ Full chain production and analysis of event data
 - ◆ Prototyping of typical physicist analysis sessions
 - ◆ Graphical monitoring display of event throughput throughout the Grid
 - ◆ Live update display of distributed histogram population from Athena
 - ◆ Virtual data re-materialization from Athena
 - ◆ Grappa job submission and monitoring

Year 2 Plan: SDSS

- Challenge Problem 1 – Balanced resources
 - ◆ Cluster Galaxy Cataloging
 - ◆ Exercises virtual data derivation tracking
- Challenge Problem 2 – Compute Intensive
 - ◆ Spatial Correlation Functions and Power Spectra
 - ◆ Provides a research base for scientific knowledge search-engine problems
- Challenge Problem 3 – Storage Intensive
 - ◆ Weak Lensing
 - ◆ Provides challenging testbed for advanced request planning algorithms

Relationships to Other Projects



Carl Kesselman
University of Southern California
carl@isi.edu

GriPhyN External Advisory Committee Meeting
Gainesville, Florida
Jan. 7, 2002

Types of Relationships

- Technology Providers (TP)
 - ◆ Developing technology that feeds into GriPhyN
- Technology Consumers (TC)
 - ◆ Relying on results of the GriPhyN project (VDT).
- Testbeds and Infrastructure (TI)
 - ◆ Provide evaluation/execution environment for GriPhyN experiments
- Application Partners (AP)



Goals of Coordination

- Minimize potential for replication of effort
- Maximize impact of GriPhyN research
- Expand the scope of research that we can conduct

U.S. GRIDS Center (TP)

- **GRIDS: Grid Research, Integration, Deployment, & Support**
 - ◆ www.grids-center.org
- **Part of NSF Middleware Initiative**
 - ◆ Integration participants are GRIDS Center and Internet-2
- **NSF-funded center to provide**
 - ◆ State-of-the-art middleware infrastructure to support national-scale collaborative science and engineering
 - ◆ Integration platform for experimental middleware technologies
- **ISI, NCSA, SDSC, UC, UW + commercial partners**
- **NMI Software releases provide foundation for VDT**
 - ◆ GriPhyN early target for NMI software
- **Mature VDT elements will be folded into NMI releases**
- **NMI initiating cross agency Grid coordination**

Earth Systems Grid (TP)

- Data Grid focused on global change applications
- Funded under DoE SciDAC program
- Technology development
 - ◆ Data movement
 - ◆ Replica management
 - ◆ Community authorization services
- Replica and security technologies being used by GriPhyN

UK eScience Initiative (TC)

- Broad-based Grid initiative across range of application domains
- Shared infrastructure base with GriPhyN
 - ◆ Globus, SRB, Condor
- UK eScience fellows to GriPhyn/iVDGL project
 - ◆ 6 funded every year, depending on success of search
 - ◆ Advertisement out this week
- Foster/Kesselman serve on Technical Advisory Group (TAG)
- Avery on PPARC Grid Steering Committee

Particle Physics Data Grid (AP)

- Focus on application of basic data grid mechanisms to a range of HEP experiments
 - ◆ Includes applications not part of GriPhyN (BABAR, D0, JLAB)
- Funded under DoE SciDAC program
- Overlap of participants with GriPhyN
 - ◆ Both applications and CS
- Shared architecture with GriPhyN
 - ◆ Participating in DGA



National Virtual Observatory (AP,TC)

- Funded under last round of ITR proposals
- Focus on catalog federation
 - ◆ Data models, catalog mediation, etc.
- GriPhyN participants in NVO
 - ◆ SDSS application
 - ◆ Middleware and system architecture
- Common Grid infrastructure
- Potential consumer of GriPhyN technology

"Infrastructure" Data Grid Projects

→ GriPhyN (US, NSF)

- ◆ Petascale Virtual-Data Grids
- ◆ <http://www.griphyn.org/>

→ Particle Physics Data Grid (US, DOE)

- ◆ Data Grid applications for HENP
- ◆ <http://www.ppdg.net/>

→ European Data Grid (EC, EU)

- ◆ Data Grid technologies, EU deployment
- ◆ <http://www.eu-datagrid.org/>

→ TeraGrid Project (US, NSF)

- ◆ Dist. supercomp. resources (13 TFlops)
- ◆ <http://www.teragrid.org/>

→ iVDGL + DataTAG (NSF, EC, others)

- ◆ Global Grid lab & transatlantic network

- ◆ Collaborations of application scientists & computer scientists
- ◆ Focus on infrastructure development & deployment
- ◆ Broad application

European Data Grid (TI)

- Complementary to GriPhyN
 - ◆ Focus on integration and applications, not research
 - ◆ Element of newly announced LHC Grid
- Initial DataGrid testbed constructed
 - ◆ Based on Globus V2.0
- Potential consumer of GriPhyN technologies
- Large overlap in application communities
 - ◆ CMS, ATLAS
- Active collaboration with GriPhyN CS project members
 - ◆ E.g. replica management
 - ◆ Foster and Kesselman serve on EDG management board

TeraGrid (TI)

- National scale infrastructure for computational science
 - ◆ Joint project between NPACI and The Alliance (NCSA)
- Next generation computational environment with focus on Grid-based applications
- Heterogeneous resource base
 - ◆ Computer, storage, graphics
- GriPhyN called out as target TeraGrid application community
- Common Grid infrastructure elements between TeraGrid & GriPhyN
 - ◆ NMI, SRB, etc.

iVDGL Summary Information

→ GriPhyN + PPDG project

- ◆ NSF ITR program \$13.65M + \$2M (matching)

→ Principal components (as seen by USA)

- ◆ Tier1 + proto-Tier2 + selected Tier3 sites
- ◆ Fast networks: US, Europe, transatlantic (DataTAG), transpacific?
- ◆ Grid Operations Center (GOC)
- ◆ Computer Science support teams
- ◆ Coordination with other Data Grid projects

→ Experiments

- ◆ HEP: ATLAS, CMS + (ALICE, CMS Heavy Ion, BTeV, others?)
- ◆ Non-HEP: LIGO, SDSS, NVO, biology (small)

→ Proposed international participants

- ◆ 6 Fellows funded by UK for 5 years, work in US
- ◆ US, UK, EU, Japan, Australia (discussions with others)

HEP Grid Coordination Effort (HICB)

→ Participants in HICB

- ◆ GriPhyN, PPDG, iVDGL, TeraGrid, EU-DataGrid, CERN
- ◆ National efforts (USA, France, Italy, UK, NL, Japan, ...)

→ Have agreed to collaborate, develop joint infrastructure

- ◆ 1st meeting Mar. 2001 Amsterdam (GGF1)
- ◆ 2nd meeting Jun. 2001 Rome (GGF2)
- ◆ 3rd meeting Oct. 2001 Rome
- ◆ 4th meeting Feb. 2002 Toronto (GGF4)

→ Coordination details

- ◆ Joint management, technical boards, open software agreement
- ◆ Inter-project dependencies, mostly High energy physics
- ◆ Grid middleware development & integration into applications
- ◆ Major Grid and network testbeds ⇒ iVDGL + DataTAG

Complementary

Global Grid Forum

- Promote Grid technologies via "best practices," implementation guidelines, and standards
- Meetings three times a year
 - ◆ International participation, hundreds of attendees
- Many GriPhyN participants contributing to GGF
 - ◆ Working group chairs, document production, etc.
- Mature GriPhyN technologies should transition to GGF

Education and Outreach Activities of GriPhyN

Manuela Campanelli, Joe Romano
University of Texas at Brownsville

- GriPhyN E/O team:

- » Manuela Campanelli (E/O coordinator)
- » Joe Romano
- » All GriPhyN members!

- iVDGL E/O team:

- » Manuela Campanelli
- » Keith Baker
- » Tim Olson
- » Rick Cavanaugh

- UT Brownsville:

- » MSI (90% of students are Hispanic)
- » Close ties with LIGO (LSC member institution)
- » CREST proposal

- iVDGL Tier3 Centers:

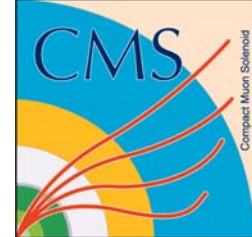
- » UT Brownsville
- » Hampton University
- » Salish Kootenai College

- **Hire E/O coordinator for GriPhyN:**
 - » UT Brownsville administration approved the hiring of a full-time tenure-track faculty member to serve as GriPhyN E/O coordinator.
 - » Search began in Fall 2000; search complete in Spring 2001.
 - » Manuela Campanelli hired to start Fall 2001:
 - Attended GGF1 and GriPhyN All-Hands meetings.
 - Made contacts with heads of dissemination programs for European data grid projects (e.g., DataGrid, EuroGrid) and ThinkQuest.
 - Helped prepare the iVDGL proposal, which will increase MSI participation in grid-related activities.
 - Interview in *'Financial times'* in Germany, and *'Brownsville Herald'* in Texas.
- **Construct UT Brownsville Linux cluster:**
 - » Benchmarking for Linux cluster begins in Fall 2000 (primarily for LIGO data analysis, but can also serve as a testbed for GriPhyN software).
 - » Construction of 96-node cluster complete in Fall 2001.

- **Development of E/O Web Site:**
 - 'Education and Outreach Center of the Grid Physics Network'*
 - <http://www.aei-potsdam.mpg.de/~manuela/GridWeb/main.html>
 - » Basic educational material about data grids, information about the physics experiments, etc.
 - » Site has already been visited by several journalists and students.
- **Grid-enable UT Brownsville Linux cluster:**
 - » Undergraduate student at UTB currently learning how to install and run Condor, Globus, gridftp, etc.
 - » Work with UW Milwaukee as part of LIGO-GriPhyN.
 - » Will be hiring a new graduate student in physics at UT Brownsville, who is very excited to do grid-related work!
- **Research Experience for Undergraduates (REU):**
 - » Preparing a proposal for an REU supplement from NSF to support 10 to 20 undergraduate students doing grid-related research during the summer.
 - » Physics education faculty member at UTB would also like to participate.
 - » Need to know who would like to mentor students? What projects?

- **Extend E/O web site:**
 - » Add more technical support information (e.g., documentation, users manuals, "how-to" guides) as virtual data toolkits become ready.
 - » Add web-based interface for accessing real data, illustrating some concepts of virtual data (e.g., <http://skyserver.sdss.org> website for SDSS data).
- **Increase MSI participation:**
 - » iVDGL funding means more minority serving institution involvement.
 - » Construct small clusters (Tier3 centers) at UT Brownsville, Hampton University, and Salish Kootenai College.
 - » Support additional undergraduate and graduate students at these institutions.
 - » Upgrade UT Brownsville cluster to 128 nodes.
- **Workshops and tutorials:**
 - » Organize "how-to" workshops and tutorials at Tier3 centers.
 - » "All-Hands Meeting" at a MSI to get a large number of minority students directly involved. (Propose UT Brownsville in Spring or Fall 2003.)

- **Leverage existing E/O programs:**
 - » QuarkNet centers at Indiana, U Florida, UT Arlington, Hampton U. Incorporate grid-related components into already existing activities. Keith Baker will lead this activity.
 - » EOT-PACI: Manuela Campanelli will work with Valerie Taylor (PI of Coalition to Diversify Computing project) to link E/O activities.
 - » ThinkQuest: Develop special challenge projects based on application sciences and grid technology. Provide "sandbox" CPUs, interesting data sets for ThinkQuest competitions.
 - » Collaboration with E/O and dissemination programs of other grid projects. Possible video documentary with members of European DataGrid and PPDG.
- **Course development:**
 - » Include grid concepts in the classes we are teaching.
 - » Identify interested MSIs and match with GriPhyN senior investigators who can give talks at these institutions.



CMS-GriPhyN

Caltech: Harvey Newman, Koen Holtman,
Julian Bunn, Conrad Steenberg,
Vladimir Litvin, Suresh Singh, Edwin
Soedermaji, Takako Hickey, Eric
Aslakson

Florida: Paul Avery, Richard Cavanaugh,
Dimitri Bourilkov, Jorge Rodriguez

Wisconsin: Miron Livny, Raj Rajamani

Chicago: Jens Voekler

ANL: Mike Wilde

GriPhyN EAC Review, Gainesville, Florida

7 January, 2002

Overview

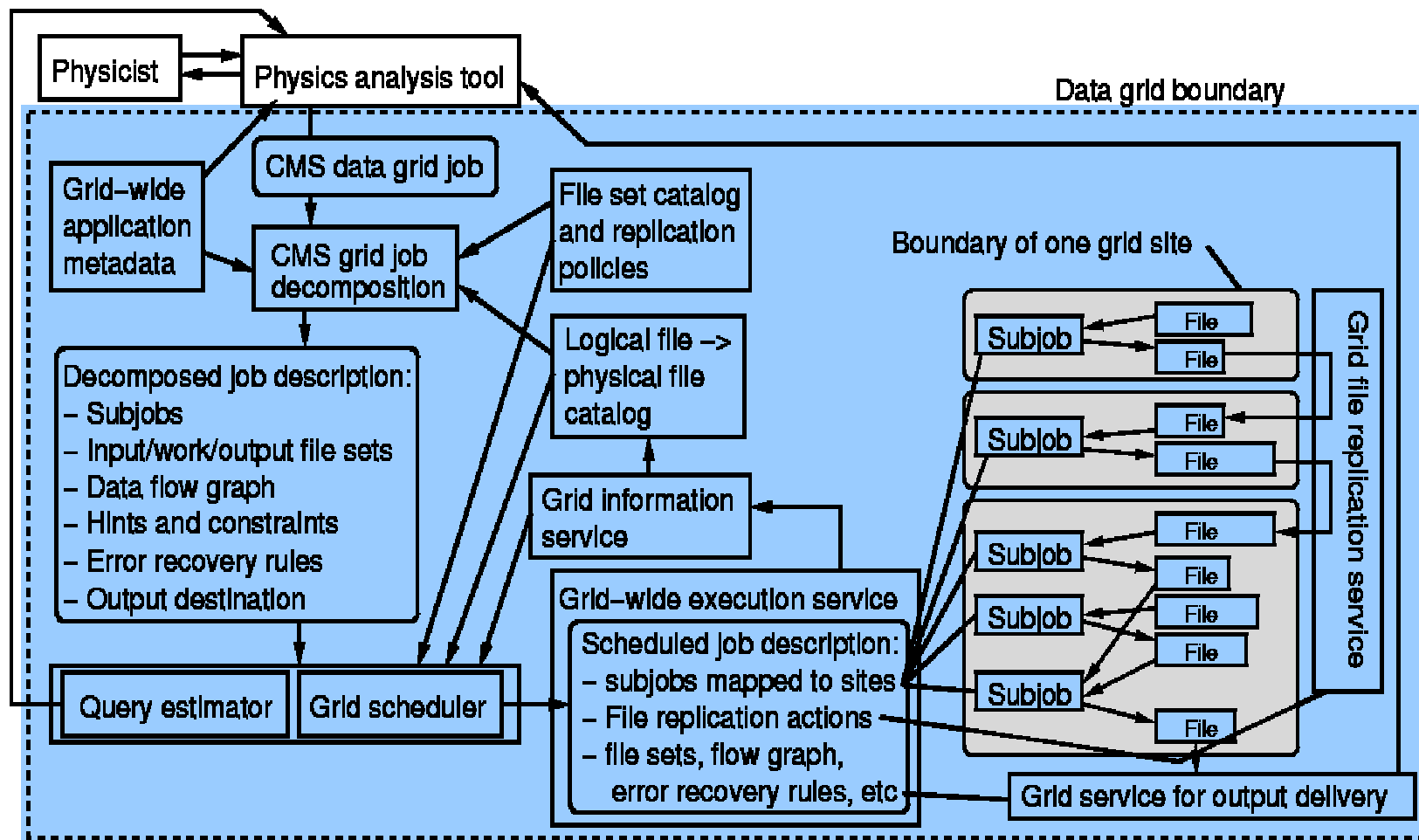
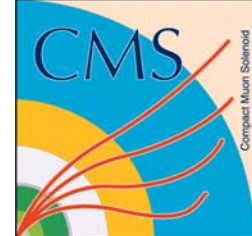
- Investigative Research (K. Holtman)
 - CMS Grid Requirements (GRIPHYN 2001-1, 2001-4)
 - Virtual Data Needs (GRIPHYN 2001-16)
- Accomplishments:
 - Demonstrated at SC2001
 - > *Virtual Data for Real Science*
 - > *Bandwidth Greedy Grid-enabled Object Collection Analysis for Particle Physics*
 - Good Collaboration with CMS Computing Teams
- Year 2 Plans
 - Simulated Data Production
 - Data Analysis
 - Test Beds



Compact Muon Solenoid



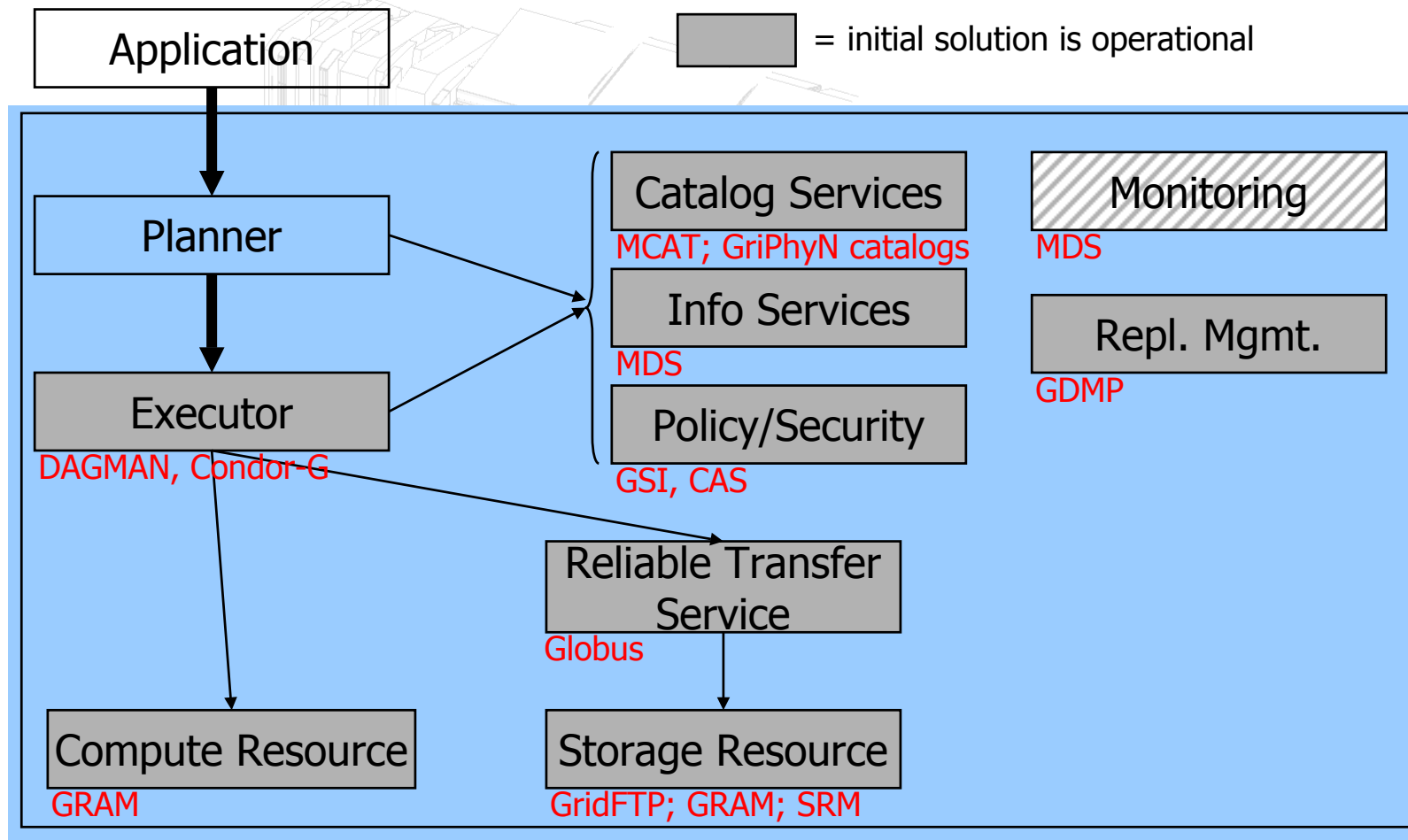
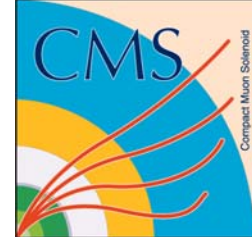
CMS Grid Requirements



Officially adopted by CMS: CMS Note 2001/037
GRIPHYN 2001-1



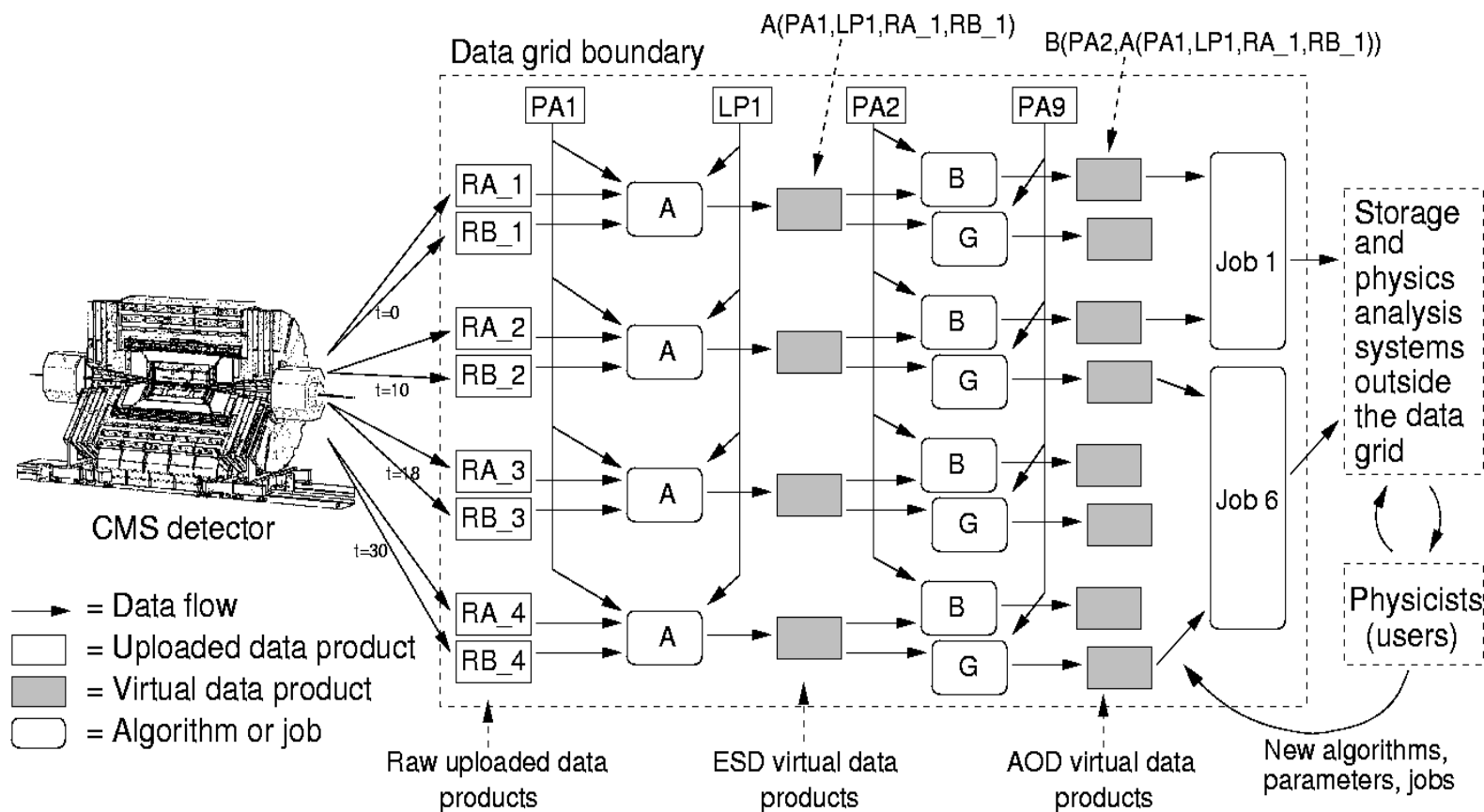
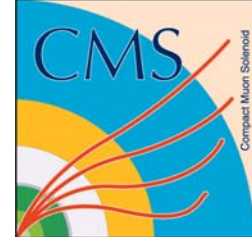
Preliminary GriPhyN Data Grid Architecture



Maps very well onto the CMS Requirements!



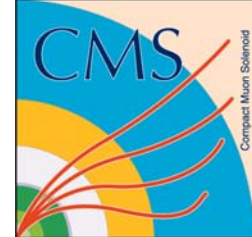
Virtual Data in CMS



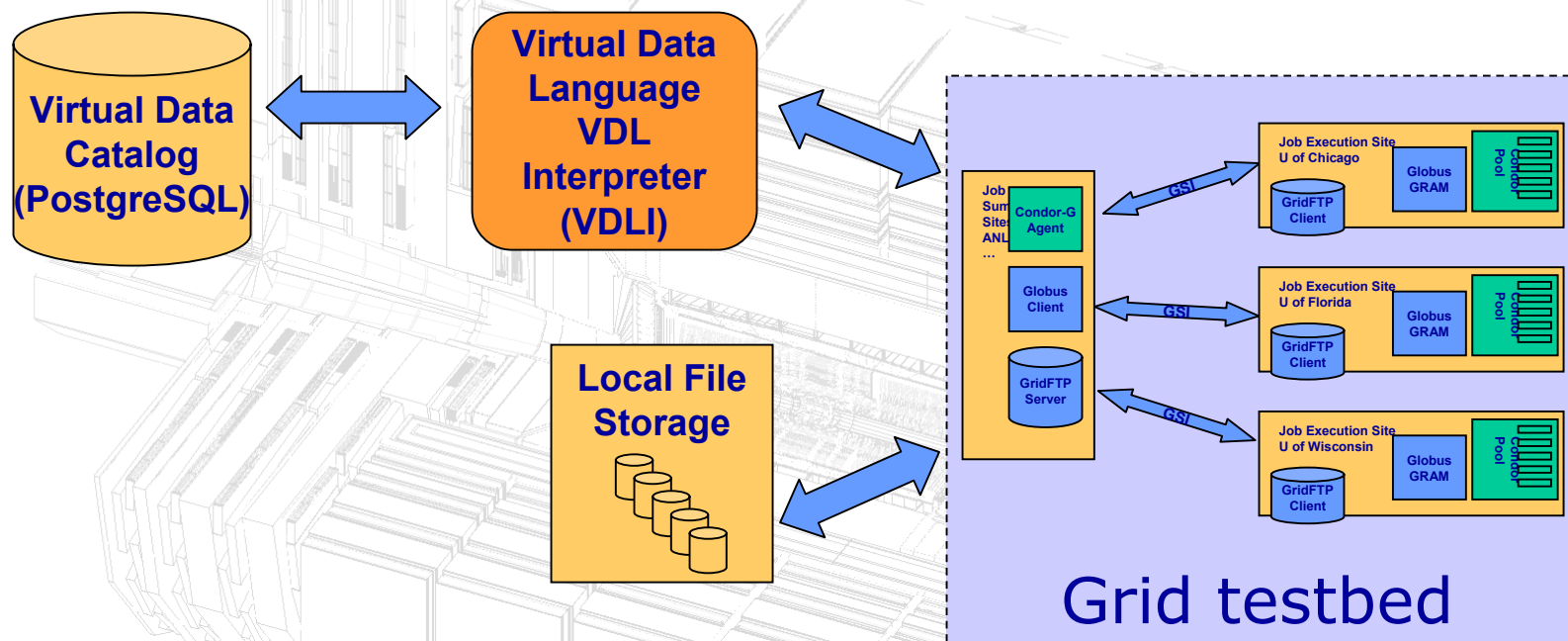
Virtual Data Long Term Vision of CMS:
CMS Note 2001/047, GRIPHYN 2001-16



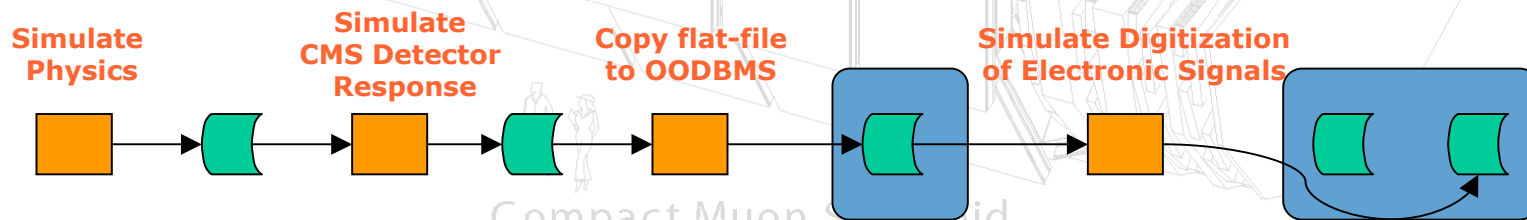
Virtual Data for Real Science: A Prototype Virtual Data Catalog



Architecture of the System:



Production DAG of Simulated CMS Data:



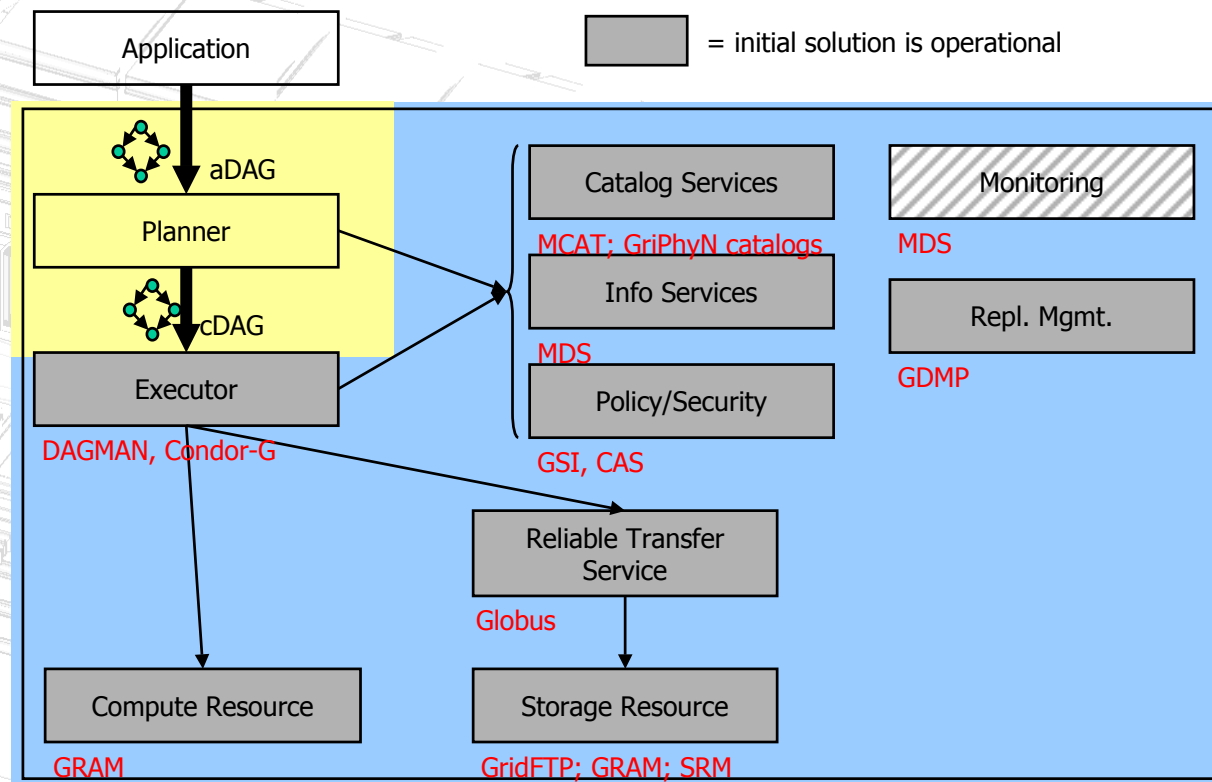
Virtual Data for Real Science: A Prototype Virtual Data Catalog

- **Abstract DAGs**

- Resource locations unspecified
- File names are logical
- Data destinations unspecified

- **Concrete DAGs**

- Resource locations determined
 - Physical file names specified
 - Data delivered to and returned from physical locations
- Translation is the job of the “planner”



Bandwidth Greedy Grid-enabled Object Collection Analysis for Particle Physics

- Physics Analysis Using a Virtual Data API

- central to the CMS experiment computing model

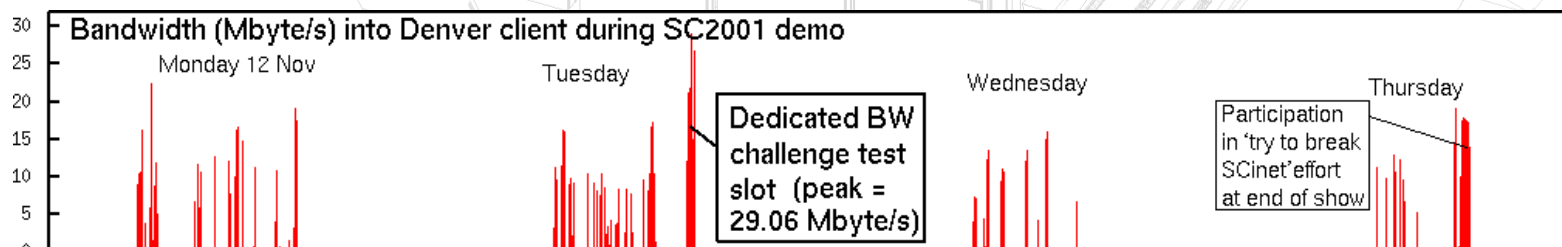
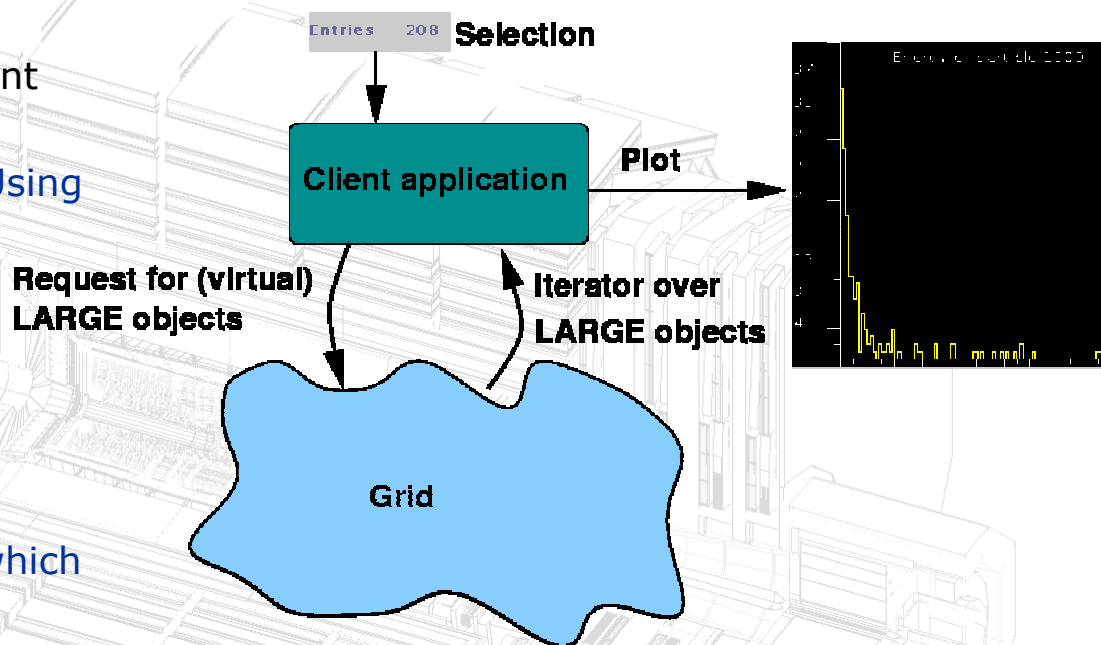
- Implemented a Virtual Data API Using Grid Technology

- Globus GSI FTP
- GDMP

- Integrated Grid Technology with OODBMS Technology

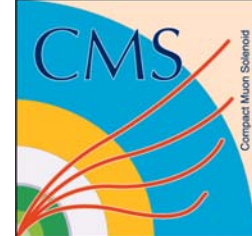
- Seamless, Easy Access to Data, which needed to:

- give all CMS physicists equal access
- successfully exploit all distributed computing resources





Production of Simulated CMS Data



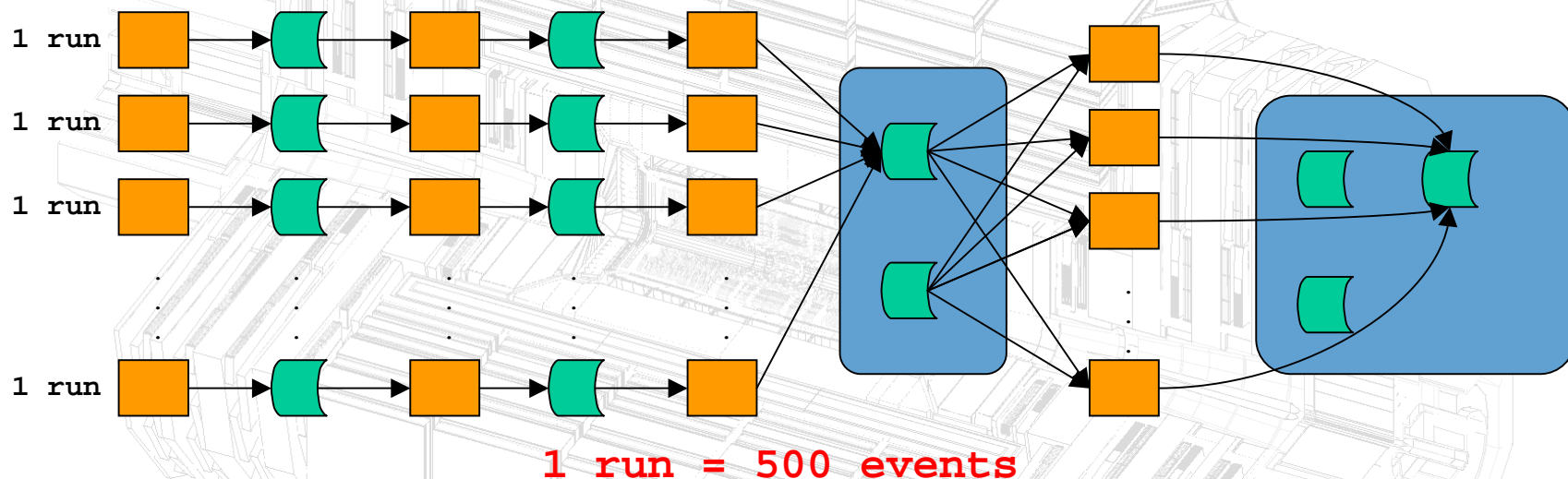
Simulate
Physics

Simulate
CMS Detector
Response

Copy flat-file
to OODBMS

Simulate Digitization
of Electronic Signals;
add noise

CPU: 2 min 8 hours 5 min 45 min



Data: 0.5 MB 200 MB 300 MB 500 MB

- IMPALA/BOSS (developed by CMS)**

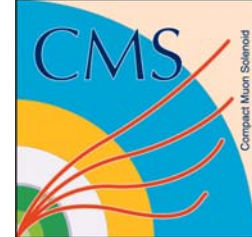
- Set of scripts for mass production of sim. data
- Provides parameter control and job tracking
- Works quite well; produced > 20 million events!
- Does not employ virtual data

- MOP (developed by PPDG)**

- Submits jobs to a Grid using DAGMan, Condor-G, Globus, and GDMP



Production of CMS Simulated Data



- CMS production of simulated data
 - $O(10)$ sites
 - $O(10^3)$ CPUs
 - ~ 20 TB of data
 - ~ 10 production managers
- Goal is to double every year—without increasing the number of production managers!
- CMS is already in an advanced stage of software development
- The immediate CMS needs are for more automation and fault tolerance



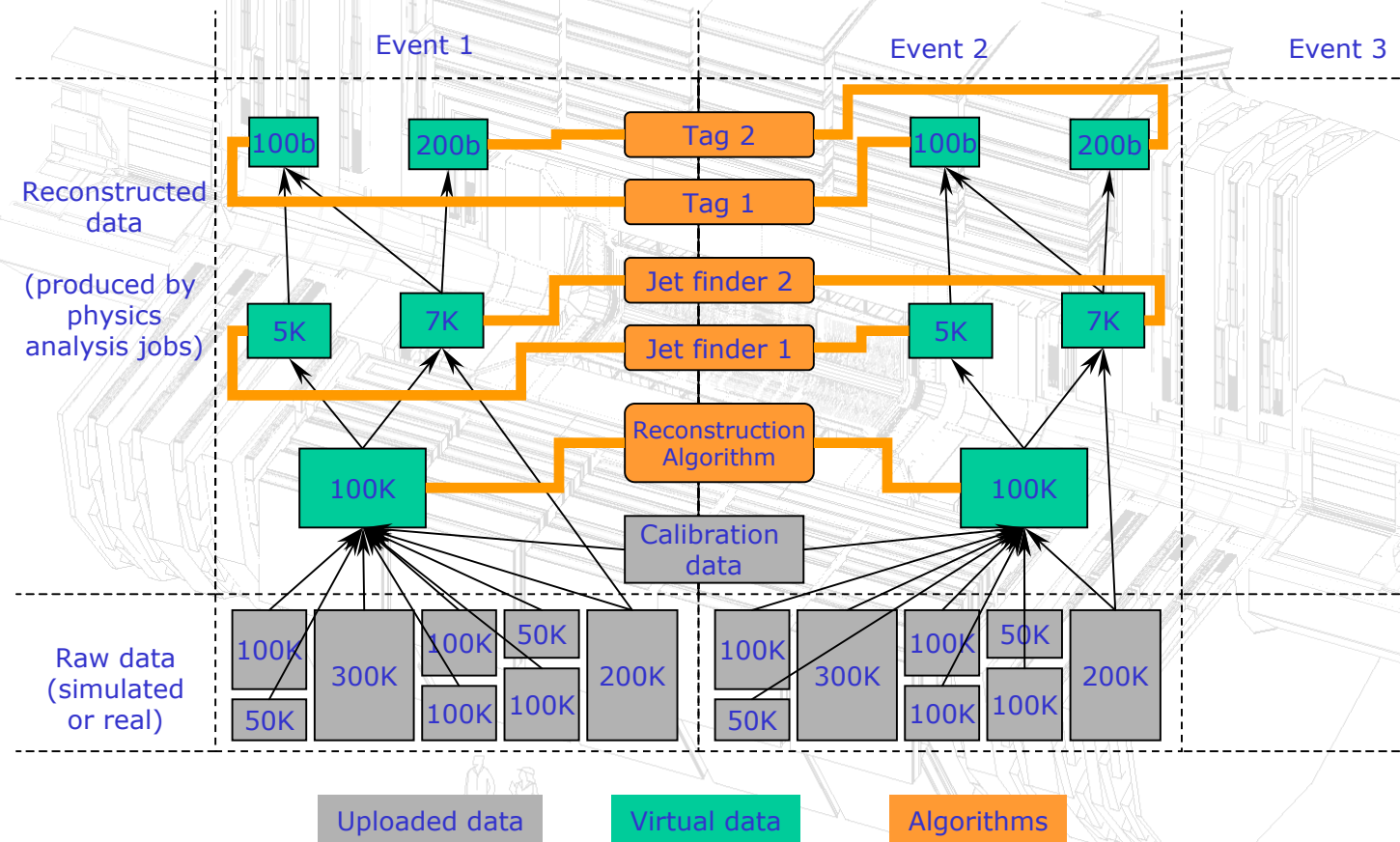
Compact Muon Solenoid

Production of Simulated CMS Data

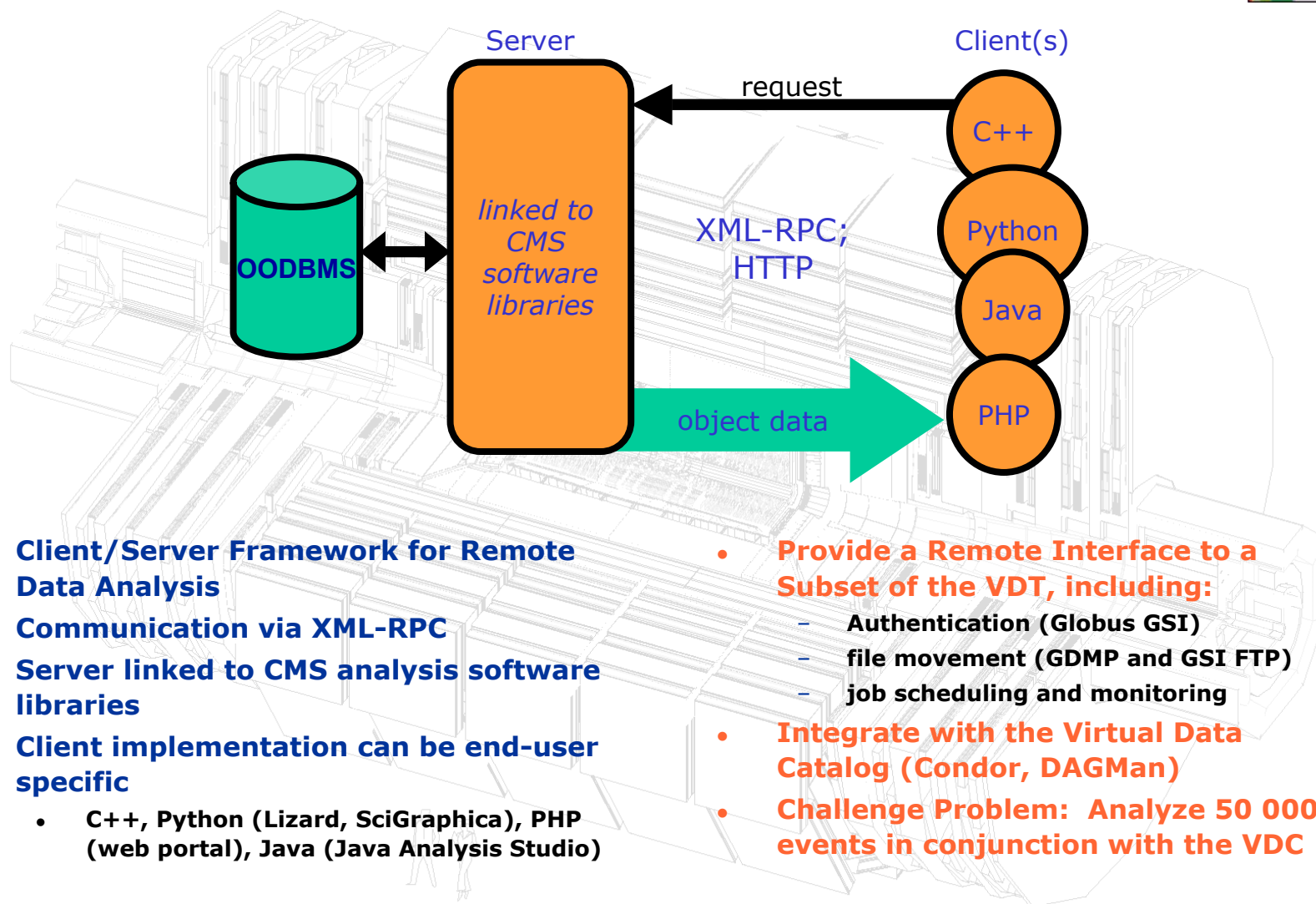
- Demonstrate the value of Virtual Data to CMS via:
 - Automatic error recovery (DAGMan)
 - Automatic data validation (add. research)
 - Automatic bookkeeping and job tracking (DB and DAGMan)
 - Transparent production at distributed sites (MOP)
- Continue to synchronize GriPhyN with CMS
 - Integrate Virtual Data into IMPALA/BOSS using an iterative, versioning process. (**bottom-up**)
 - Continue architectural work on VDC (**top-down**)
- Challenge Problems: valuable feedback for future prototyping
 - Generate 50 000 events using the VDC
 - Fulfill CERN request for simulated data production

CMS Data Analysis

Dominant use of Virtual Data in the Future



CMS Remote Data Analysis Prototype API: Clarens



- **Client/Server Framework for Remote Data Analysis**
- **Communication via XML-RPC**
- **Server linked to CMS analysis software libraries**
- **Client implementation can be end-user specific**
 - C++, Python (Lizard, SciGraphica), PHP (web portal), Java (Java Analysis Studio)
- **Provide a Remote Interface to a Subset of the VDT, including:**
 - Authentication (Globus GSI)
 - file movement (GDMP and GSI FTP)
 - job scheduling and monitoring
- **Integrate with the Virtual Data Catalog (Condor, DAGMan)**
- **Challenge Problem: Analyze 50 000 events in conjunction with the VDC**

US-CMS Test Grid



- VDT 1.0
 - Condor 6.3.1, DAGMan
 - ClassAds 0.9
 - Globus 2.0 beta
 - GDMP 2.0
- Condor-G 6.3.1
- Objectivity 6.1
- DAR (CMS Distribution After Release tarball)
 - CMS executables
 - Dynamically linked shared object libraries
 - Required parameters for simulated data production

US-CMS Test Grid

- Joint endeavor between PPDG and GriPhyN / iVDGL
 - Platform to develop grid-enabled tools for CMS
 - Use the ESNet Certificate and Registration Authority in April
 - Operational in January 2002
- Integrate with the CMS Data Grid
 - CERN, Caltech, Florida, FNAL, INFN Padova, INFN Bologna, IN2P3 Lyon, PPARC RAL
 - Important for compatibility between tools developed by CMS, GriPhyN, PPDG, and EDG.
- Integrate with existing and future GriPhyN /iVDGL Test Grids
 - Important for development of common solutions for common needs.
 - Tempered by the fact that CMS currently requires a very particular environment: Red Hat 6.2, Objectivity/DB (req. license), etc.

Conclusion

- CMS Grid Requirements Understood
- Successful
 - Development of a Prototype Virtual Data Catalog
 - Demonstration of Remote Data Analysis
 - Collaboration with CMS Computing teams
- Work closely with CMS to implement Virtual Data technology into CMS Production
- Continue Development of Data Analysis Tools
 - Dominant use of virtual data in the future
- Integrated US-CMS Test Grid with CMS Data Grid will help forge closer GriPhyN/PPDG and CMS ties.



Compact Muon Solenoid

Realizing LIGO Virtual Data

Caltech: Kent Blackburn, Phil Ehrens, Albert Lazzarini, Roy Williams

ISI: Ewa Deelman, Carl Kesselman, Gaurang Mehta, Leila Meshkat, Laura Pearlman

UWM: Bruce Allen, Scott Koranda

Outline

- Year 1 Accomplishments

- LIGO Virtual Data requirements (GriPhyN-2001-6)
- GriPhyN/LIGO Virtual Data prototype demonstration at SC'01 (GriPhyN-2001-18)
- Transformation Catalog Design (GriPhyN-2001-17)
- Outreach activities, Collaboration with other gravitational wave communities

- Plans for Year 2

- Pulsar search, as a science focus
- Virtual data
- Request planning and execution

- Issues and Challenges

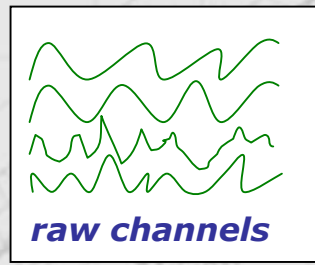
The Physics of LIGO's Pulsar Search

- Hypothesis: GW Pulsars as sources of strong GW signature
 - Presently unidentified sources
 - Support for the so-called "blind" or "all-sky" search for new sources.
- GW signals are frequency modulated by Doppler shift produced by Earth's rotation and barycentric motion around the Sun.
 - Need to account for this as a function of different lines of sight to putative sources, with different parameters.
 - Exploiting full SNR potential of LIGO data becomes a peta-flops class problem.
- Great potential for use of Virtual Data and Grid technologies
 - Many data products are "reusable"
 - Target of next phase of GriPhyN-LIGO Virtual Data research.

LIGO's Pulsar Search

(Laser Interferometer Gravitational-wave Observatory)

Interferom



Single Frame

Extract
channel



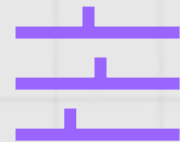
Short time frames

transpose



Long time frames

SFT



Extract
frequency
range

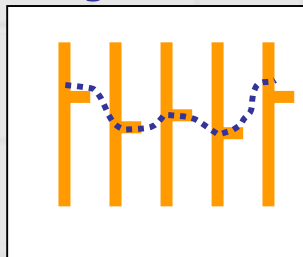


GriPhyN meeting 1/2002

Construct
image

Hz

**Time-frequency
Image**



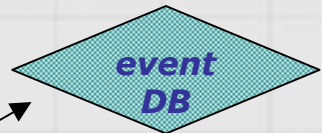
Time

Find Candidate



Store

Ewa Deelman, ISI



GriPhyN/LIGO prototype functionality

LIGO Specific
Data
Specification

XML

GriPhyN/LIGO

XML

LIGO Data
Product

- Interpret an XML-specified request
- Acquire user's proxy credentials
- Consult replica catalog to find available data
- Construct a plan to produce data not available
- Execute the plan
- Return requested data in Frame or XML format

Year 1 Virtual Data Product: Channel Extraction

Compute resources running LIGO Data Analysis System at Caltech and UWM, storage resources at ISI, UWM and Caltech



xml

**HTTP
frontend**

**MyProxy
server**

Cgi interface

**Transformation
Catalog**

Planner

Monitoring

**Replica
Catalog**

G-DAG (DAGMan)

**Executor
CondorG/
DAGMan**

Logs

GridCVS

GridFTP

GRAM/LDAS

GridFTP

GRAM

LDAS

**Storage
Resource**

**Compute
Resource**

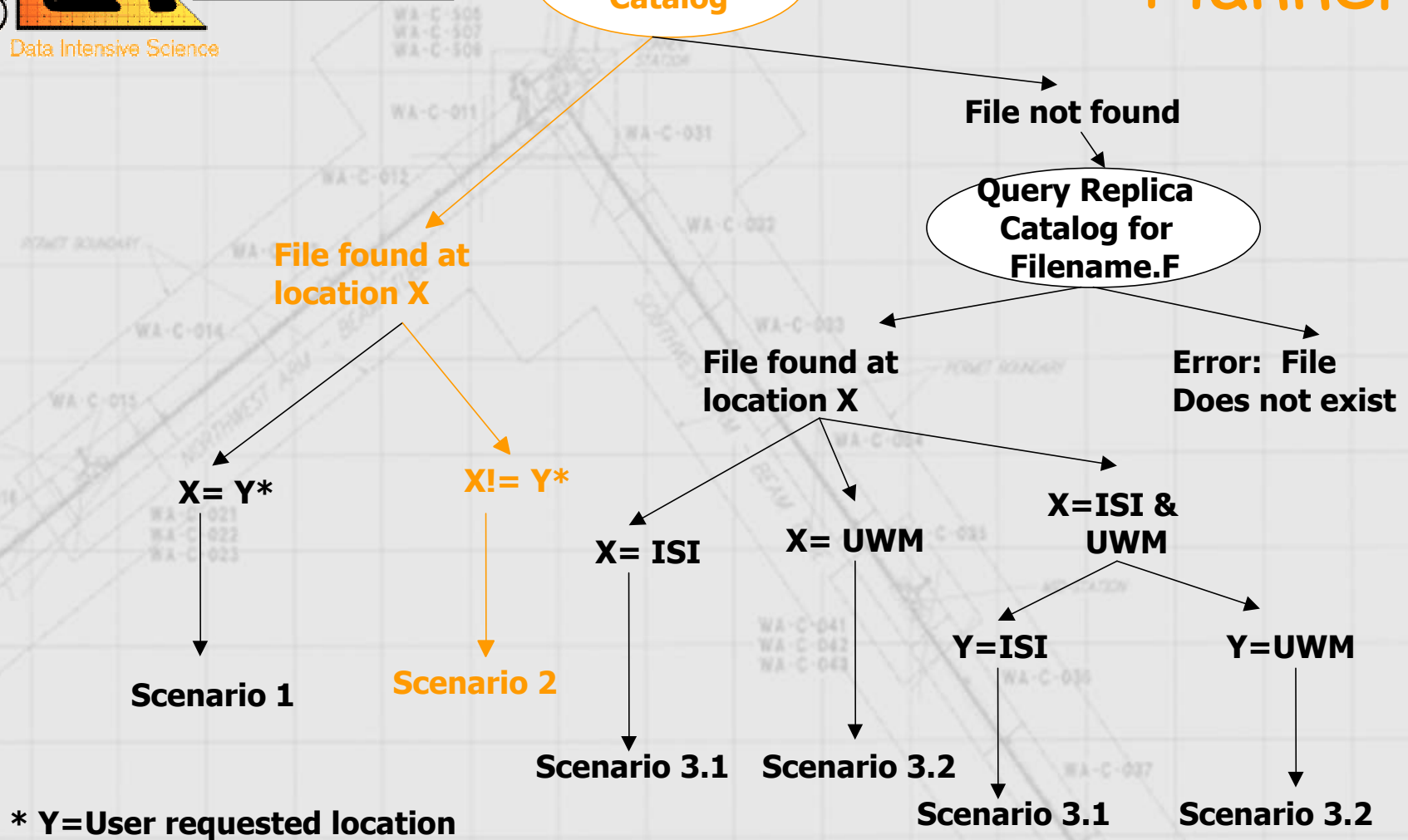
002

1

Collection name
Channel name
Filename.F
Time interval
Desired Loc

Query Replica
Catalog

Planner



* Y=User requested location

Template instantiation (scenario 2)

Information gathered by Planner:
C_A_100 in *dc.isi.edu/frames*
Output location:
host.uwm.edu/myframes

globus_url_copy X
From *a* to *b*

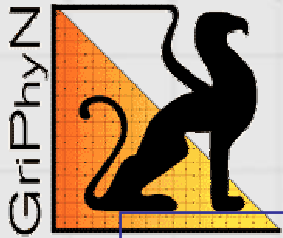
Register *X*
In RC with
location *b*

Abstract G-DAG

**Concrete G-DAG
(DAGMan)**

globus_url_copy C_A_100
From *dc.isi.edu/frames* to
To *host.uwm.edu/myframes*

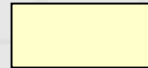
Register *C_A_100*
In RC with
location *host.uwm.edu/myframes*



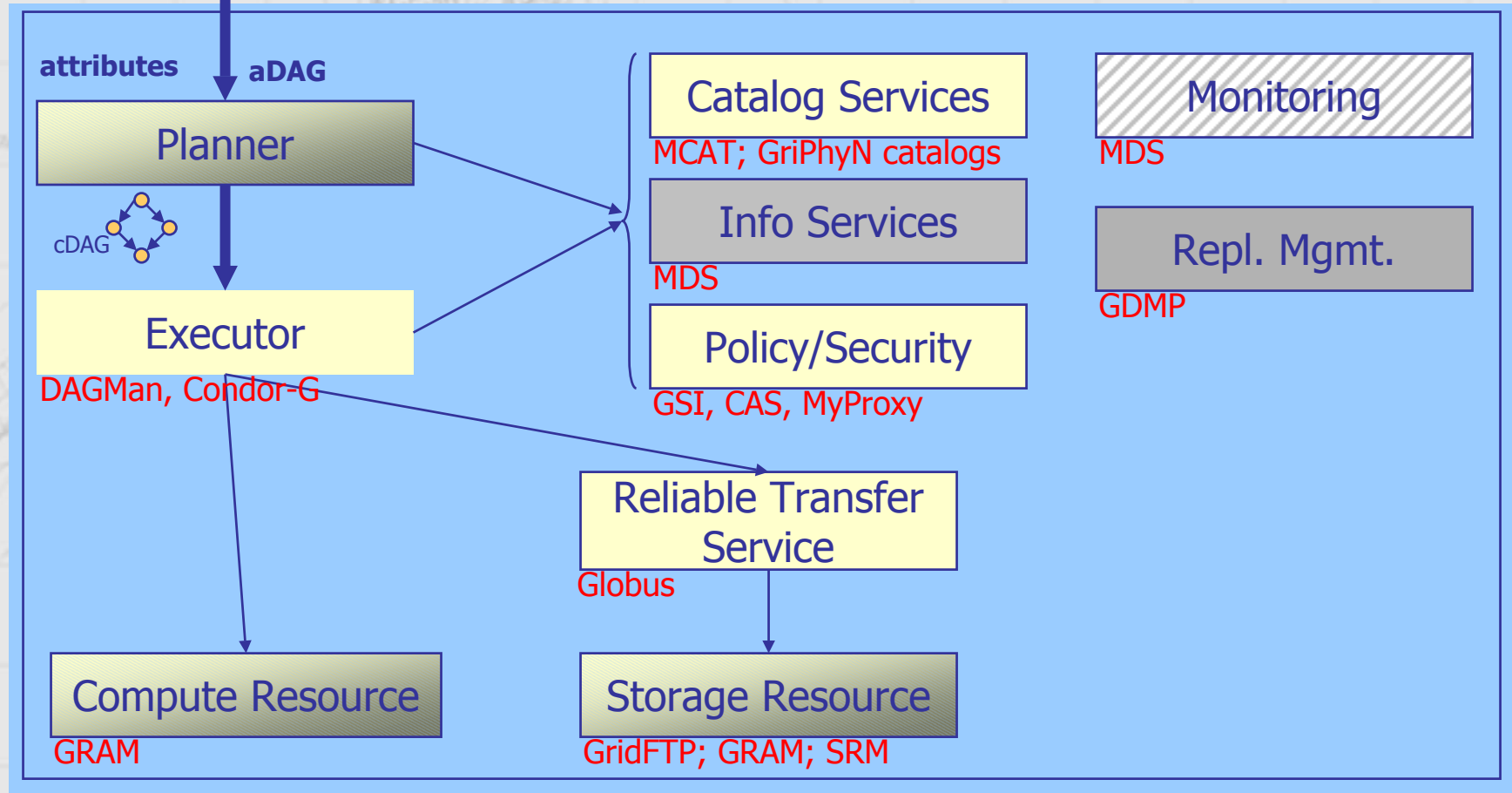
Preliminary GriPhyN Data Grid Architecture



New/modified in Prototype



Standard Globus or Condor-G component



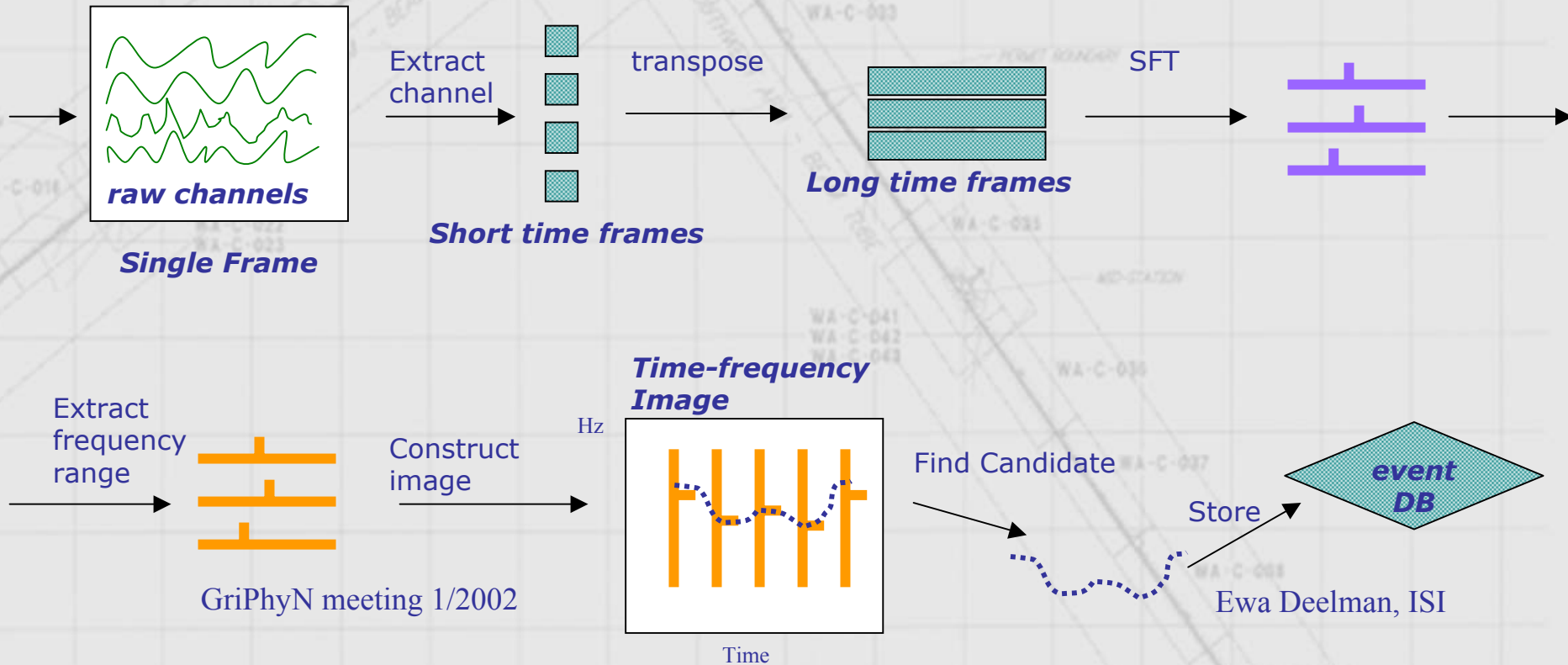
Accomplishments

- Simple demonstration of Virtual Data Concepts
 - Transparency with respect to location
 - Transparency with respect to materialization
- Provided a Globus interface to LDAS
 - Basis for a secure access to LIGO resources
- Designed the Transformation Catalog
 - Used to find an appropriate executable for a given architecture
 - Can be used in many systems
- Basic infrastructure for the development of Virtual Data concepts
 - Foundation for Year 2

LIGO's Pulsar Search

(Laser Interferometer Gravitational-wave Observatory)

Interferom



Year 2 The Year of the Pulsar Search Mock Data Challenge

- Broaden the GRAM/LDAS interface
 - greater variability and functionality: SFTs, concatenation, decimation and resampling.
- Design a Data Discovery mechanism for discovery of data replicas on a Grid.
 - ability to interact with the LDAS Diskcache resources
- Implementation of the Data Discovery mechanism to support the pulsar search

Virtual Data Concepts

- Implement the Transformation Catalog.
- Explore the design of the Derived Data Catalog, which specifies how Virtual Data products are materialized.
- Unify the catalog schemes used by CMS and LIGO – base it on a common VDT 2.0 release.
- Apply replication concepts by developing a real-time international mirror, and a fault-tolerance replica at UW-Mil.
- Use of Catalogs to materialize Virtual data required in the pulsar search (including Transformation Catalog).

Planning and Fault Tolerance

- Specify the planning requirements
- Evaluate the available solutions
- Prototype a more sophisticated planner
- Specify LIGO's fault tolerance requirements, extrapolate to GriPhyN in general
- Assess existing fault and failure issues within LIGO
- Assess the applicability of existing techniques

Year 2 Challenges

- Explore bulk data operations
 - Finding new available data
 - Registering data into catalogs
- Deepen the understanding of Virtual Data naming
 - How do you ask for what you want?
- Planning and Fault Tolerance
 - Need to specify model
 - Explore existing planning solutions
 - Examine fault tolerance issues at the system level
- Scalable pulsar search to scientifically interesting levels of sensitivity at SC'2002